

Quantity vs. Quality: Evaluating User Interest Profiles Using Ad Preference Managers

Muhammad Ahmad Bashir

Umar Farooq

Maryam Shahid

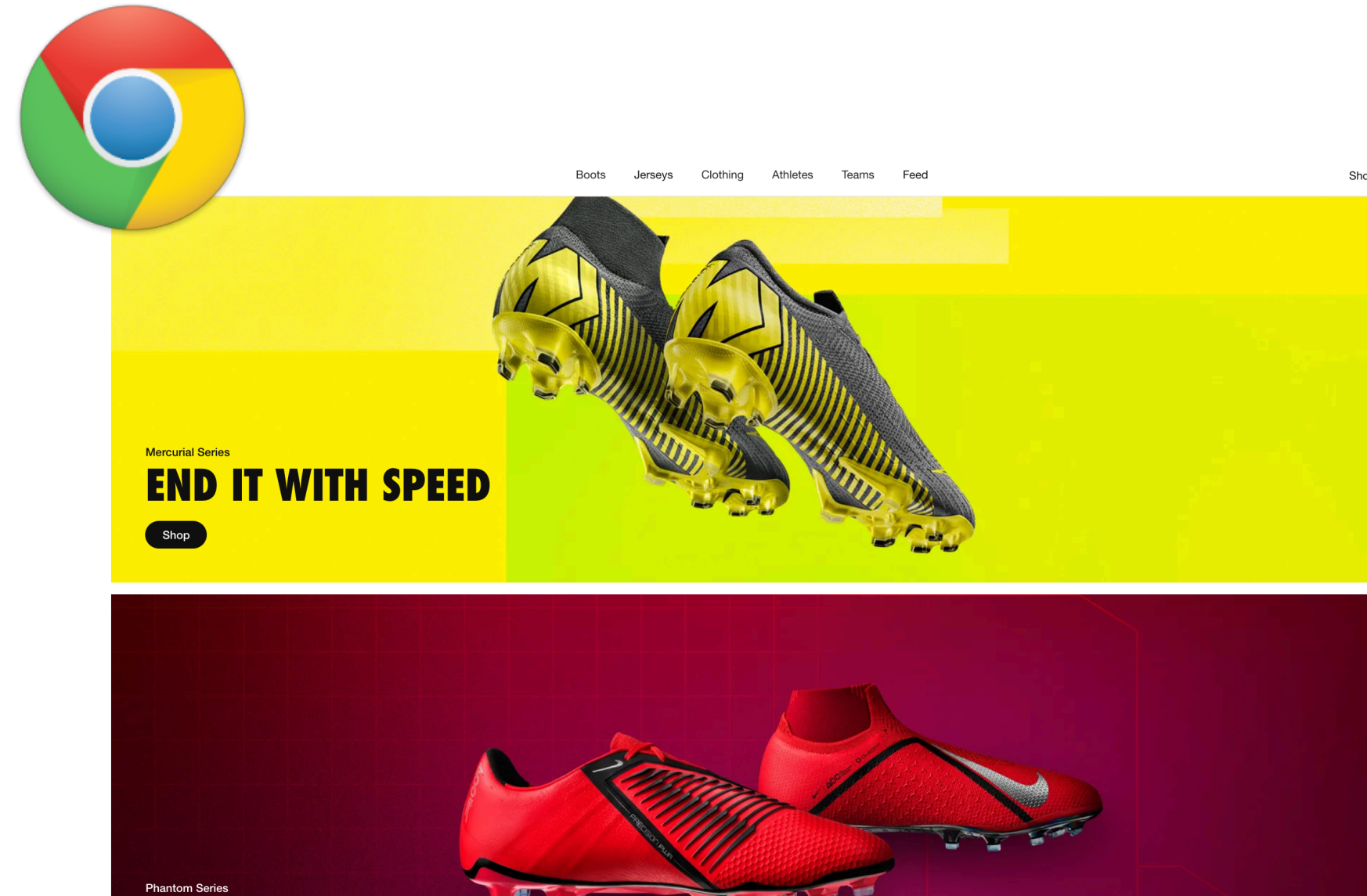
Muhammad Fareed Zaffar

Christo Wilson

Northeastern University
**Khoury College of
Computer Sciences**



Online Tracking



Online Tracking



Online Tracking



Online Tracking



Inferences Used For Targeted Ads



washingtonpost.com

Sections

The Washington Post
Democracy Dies in Darkness

Sign In

Washington Wizards

Just as the Wizards were stumbling, the new guys found their footing



Jeff Green played starter's minutes in Saturday's win. (Wilfredo Lee/Associated Press)



Most Read Sports

- 1** Redskins earn an ugly 16-3 win over the Buccaneers, remain in first place in the NFC East
- 2** **Analysis** Redskins-Buccaneers takeaways: Tampa dominates the stat sheet, including with game-altering turnovers
- 3** Saints and Chiefs roll; Baker Mayfield leads Browns to victory; Patriots slip up
- 4** **Analysis** The Patriots' path back to the Super Bowl just got more complicated
- 5** Exercise rider and horse dead after early-morning accident at Churchill Downs

- Home
- Facebook
- Twitter
- Google+
- Email
- LinkedIn
- Pinterest
- Tumblr
- Print
- Share
- 2

Inferences Used For Targeted Ads



washingtonpost.com

Sections

The Washington Post
Democracy Dies in Darkness

Sign In

Washington Wizards

Just as the Wizards were stumbling, the new guys found their footing



Jeff Green played starter's minutes in Saturday's win. (Wilfredo Lee/Associated Press)



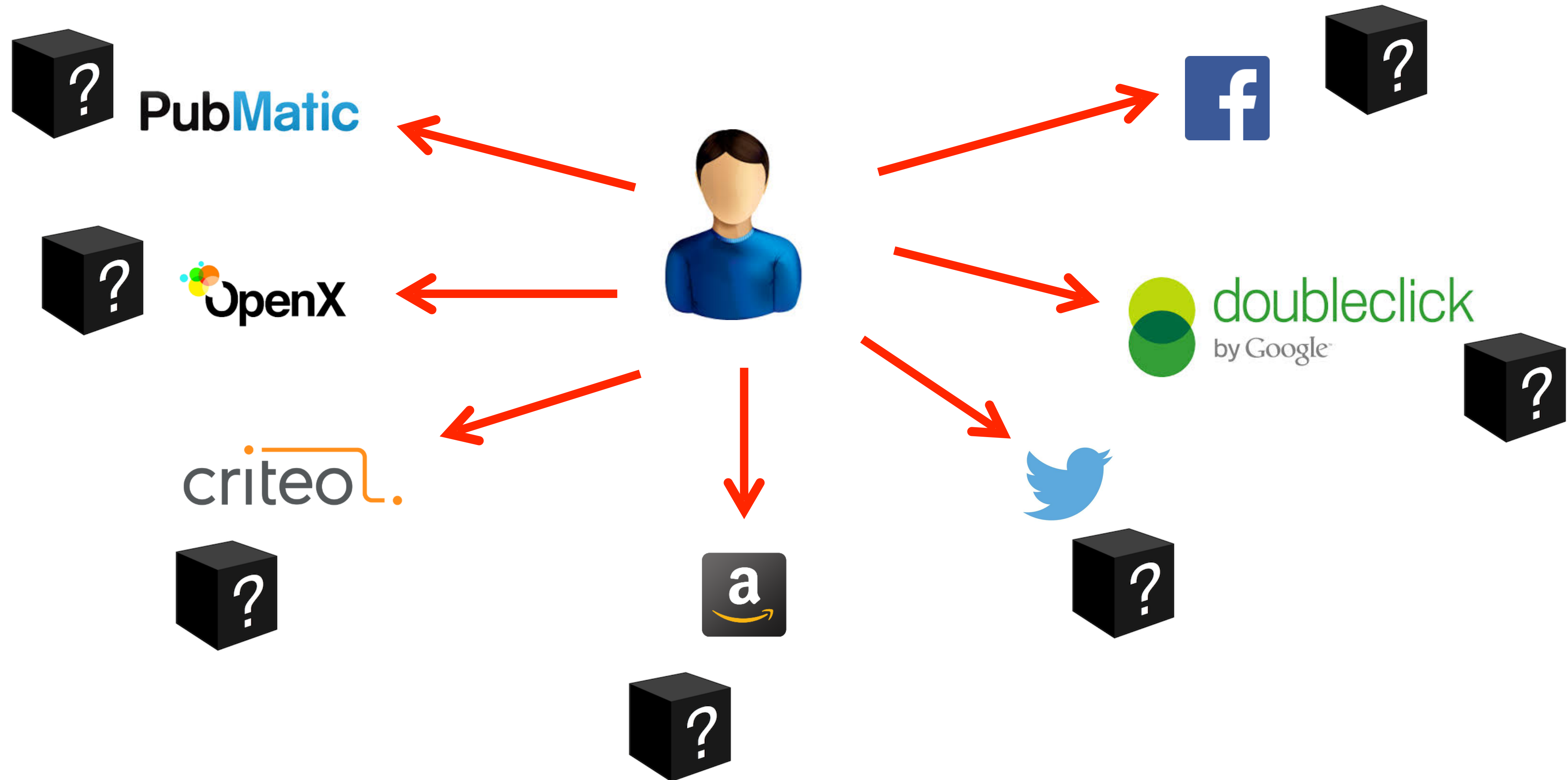
Most Read Sports

- 1** Redskins earn an ugly 16-3 win over the Buccaneers, remain in first place in the NFC East
- 2** **Analysis** Redskins-Buccaneers takeaways: Tampa dominates the stat sheet, including with game-altering turnovers
- 3** Saints and Chiefs roll; Baker Mayfield leads Browns to victory; Patriots slip up
- 4** **Analysis** The Patriots' path back to the Super Bowl just got more complicated
- 5** Exercise rider and horse dead after early-morning accident at Churchill Downs

We Don't Know What Ad Networks Infer



We Don't Know What Ad Networks Infer



Goals of the Study

1. Who knows what and how much?
2. How do users perceive interests inferred about them?
3. How are the interests inferred?
4. How do privacy practices impact amount of inferences drawn?

Goals of the Study

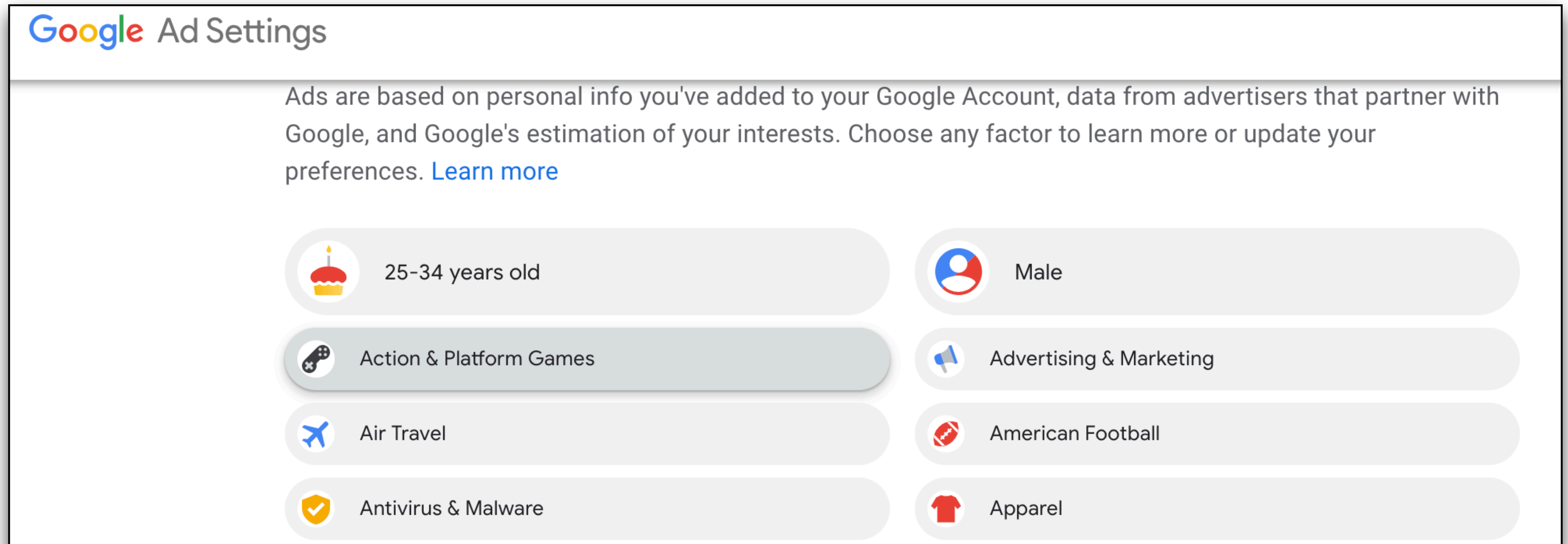
1. Who knows what and how much?
2. How do users perceive interests inferred about them?
3. ~~How are the interests inferred?~~
4. ~~How do privacy practices impact amount of inferences drawn?~~

Ad Preference Managers (APMs)

- Transparency tools
- Let users control the inferred interests about them

Ad Preference Managers (APMs)

- Transparency tools
- Let users control the inferred interests about them



The screenshot displays the 'Google Ad Settings' interface. At the top, it reads 'Google Ad Settings'. Below this, a paragraph explains: 'Ads are based on personal info you've added to your Google Account, data from advertisers that partner with Google, and Google's estimation of your interests. Choose any factor to learn more or update your preferences. [Learn more](#)'. The interface features two columns of interest categories, each with a circular icon and a text label. The categories are: '25-34 years old' (cake icon), 'Male' (person icon), 'Action & Platform Games' (game controller icon), 'Advertising & Marketing' (megaphone icon), 'Air Travel' (airplane icon), 'American Football' (football icon), and 'Antivirus & Malware' (checkmark icon).

Google Ad Settings

Ads are based on personal info you've added to your Google Account, data from advertisers that partner with Google, and Google's estimation of your interests. Choose any factor to learn more or update your preferences. [Learn more](#)

- 25-34 years old
- Male
- Action & Platform Games
- Advertising & Marketing
- Air Travel
- American Football
- Antivirus & Malware
- Apparel

Overview

1. Data collection
2. Interests inferred by different APMs
3. Perception of interests
4. Limitations & Conclusion

Data Collection

Data Collection

- We recruited 220 participants
 - 82 from Pakistan (university students), 138 from US (crowdsourcing)

Data Collection

- We recruited 220 participants
 - 82 from Pakistan (university students), 138 from US (crowdsourcing)
- Used our browser extension to
 - A. Take a survey
 - B. Contribute data from their APMs + Historical Data

Data Collection

- We recruited 220 participants
 - 82 from Pakistan (university students), 138 from US (crowdsourcing)
- Used our browser extension to
 - A. Take a survey
 - B. Contribute data from their APMs + Historical Data

Ethics

- Obtained IRB from both LUMS and Northeastern University
- Obtained informed consent.

Browser Extension

Foreground

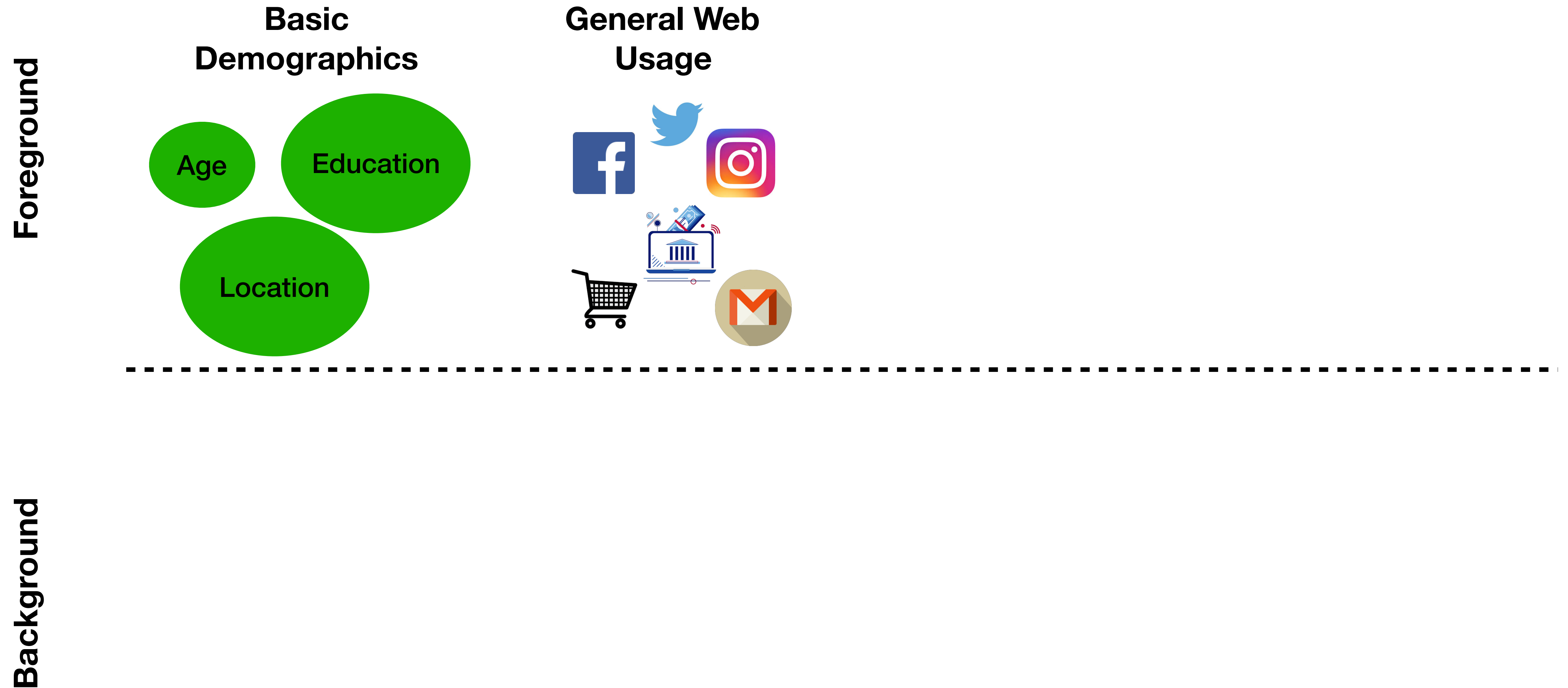
Background



Browser Extension



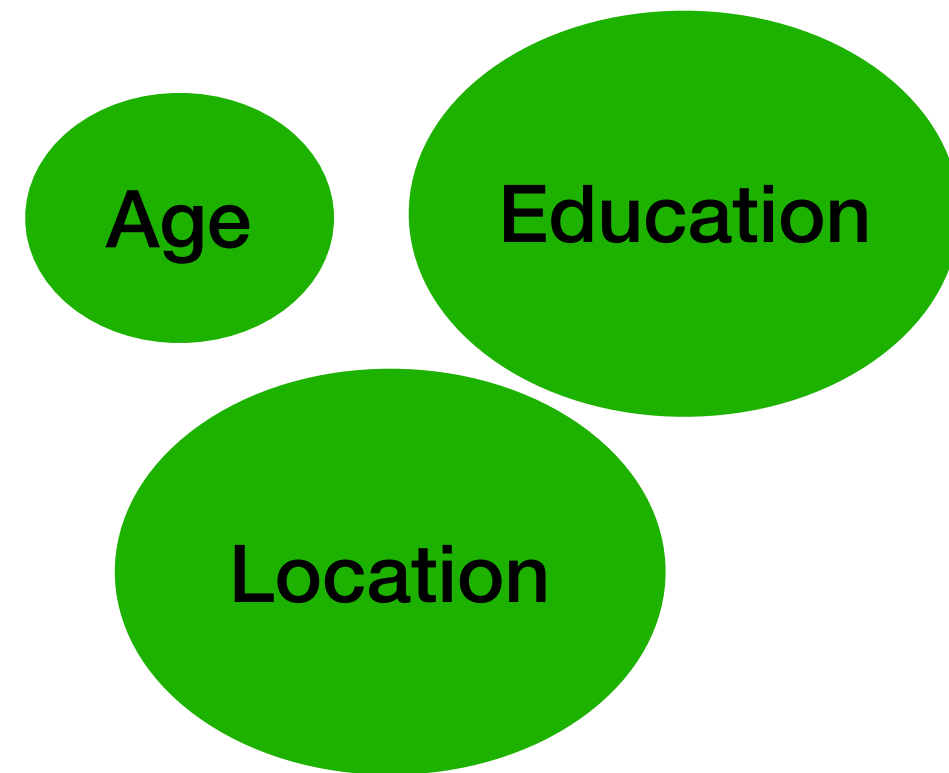
Browser Extension



Browser Extension

Foreground

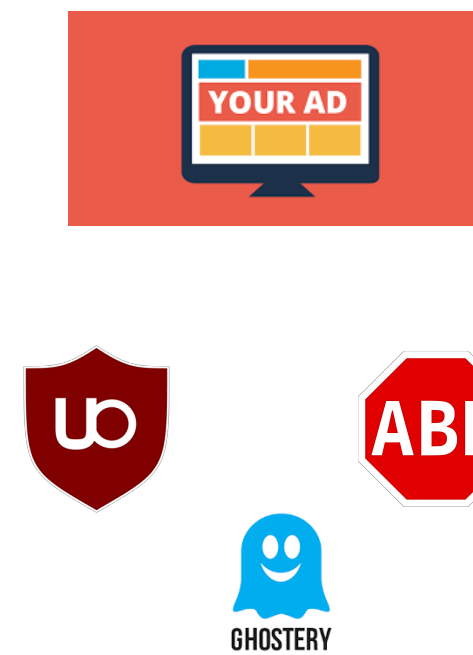
Basic Demographics



General Web Usage

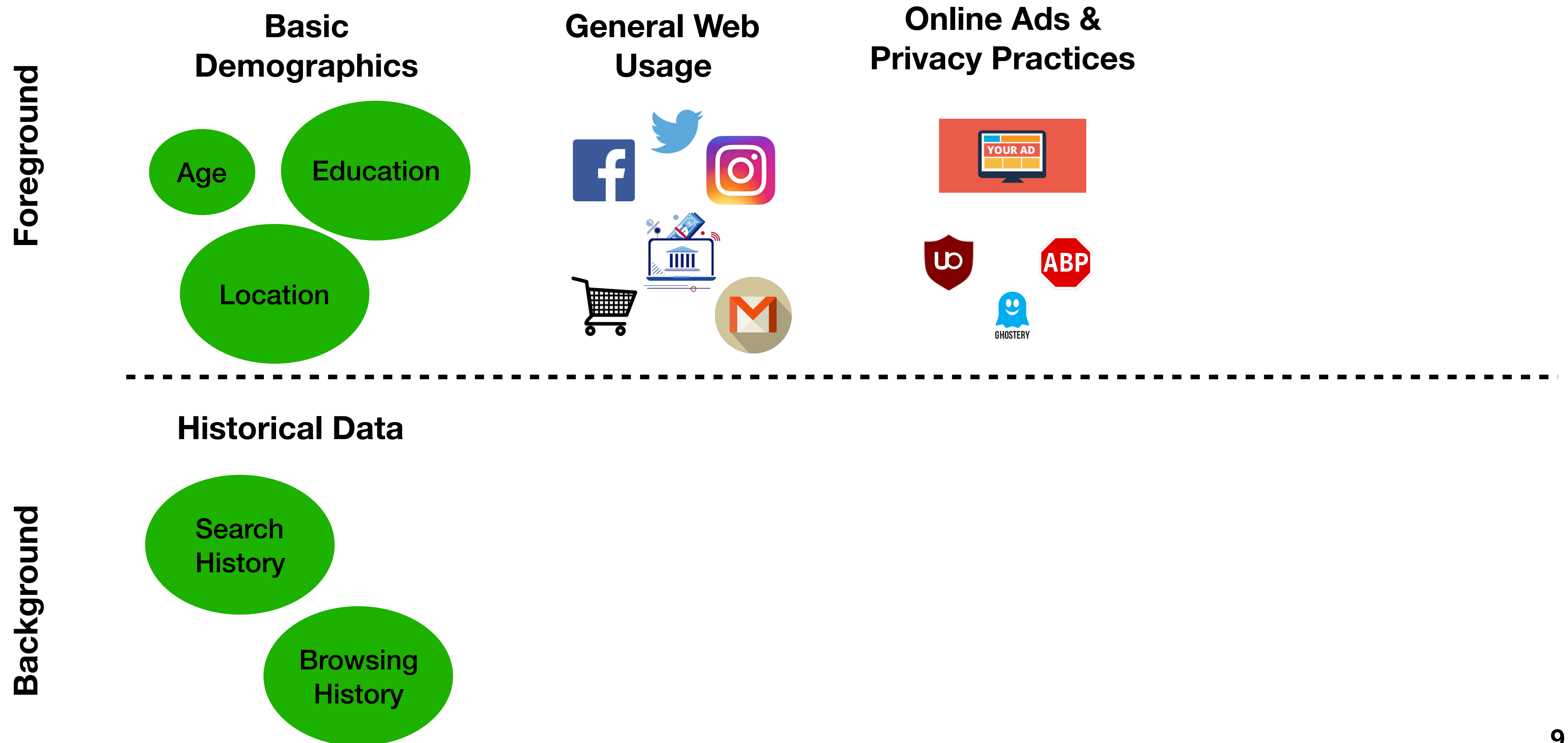


Online Ads & Privacy Practices

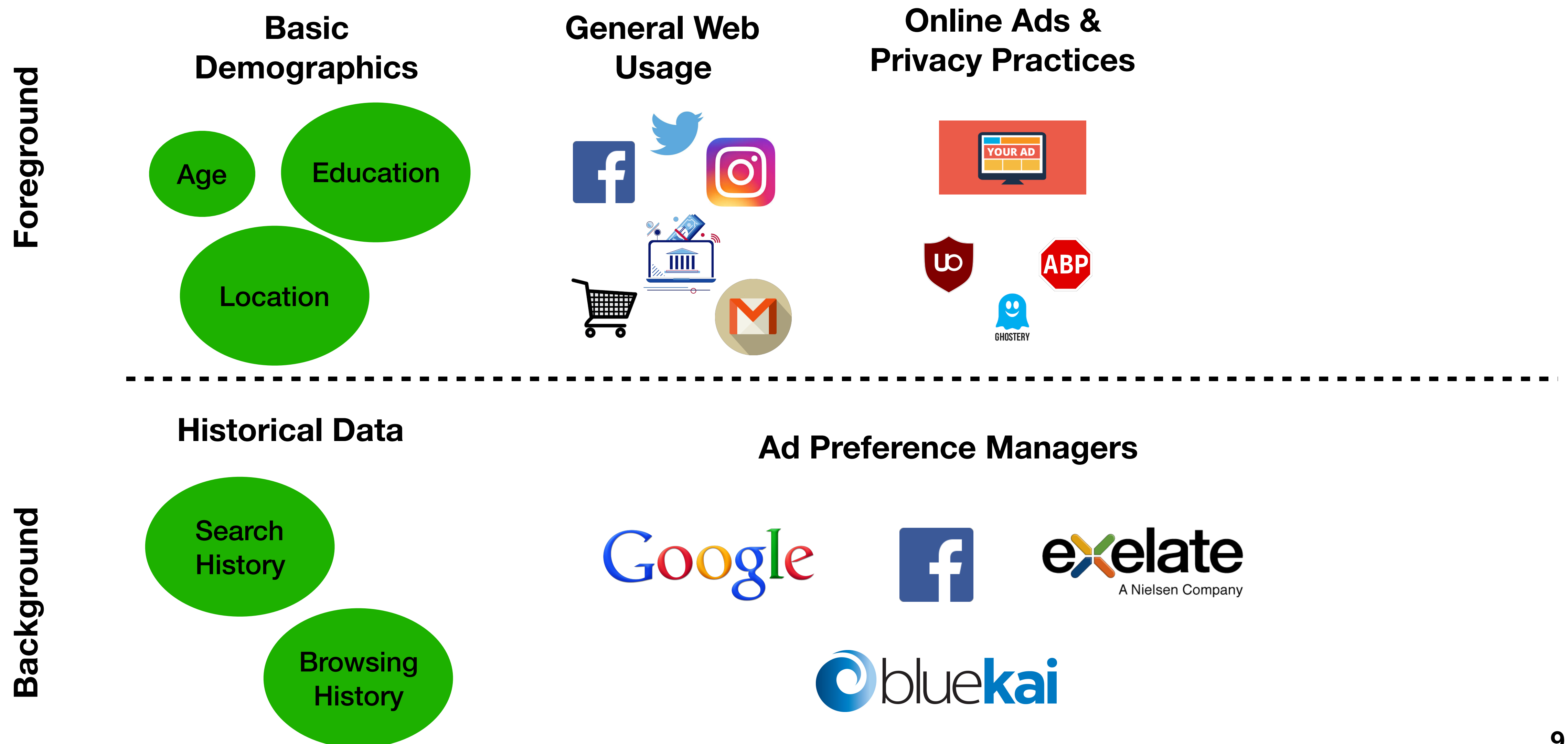


Background

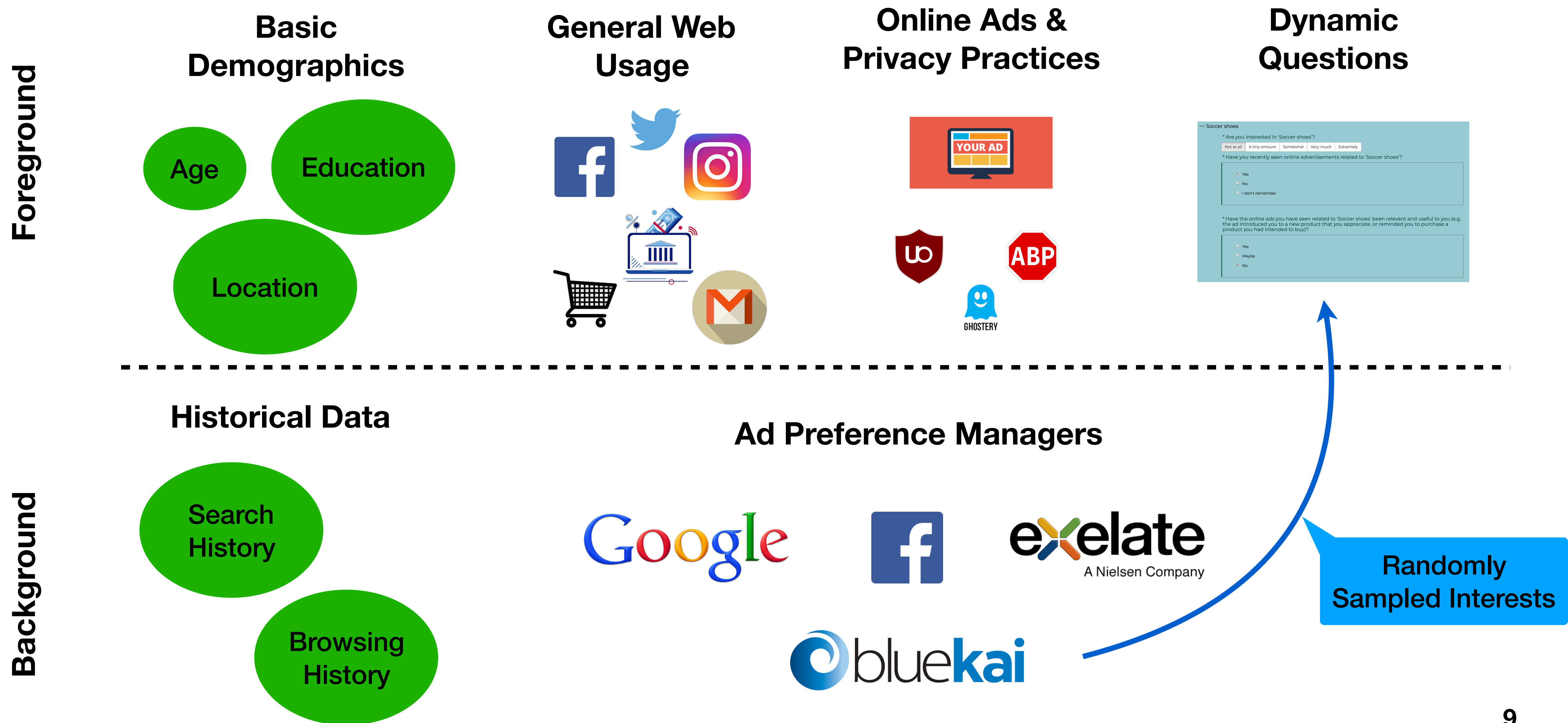
Browser Extension



Browser Extension



Browser Extension



Dynamic Questions

~~ Soccer shoes

* Are you interested in 'Soccer shoes'?

Not at all

A tiny amount

Somewhat

Very much

Extremely

* Have you recently seen online advertisements related to 'Soccer shoes'?

Yes

No

I don't remember

Dynamic Questions

~~ Soccer shoes

* Are you interested in 'Soccer shoes'?

Not at all A tiny amount Somewhat Very much Extremely

* Have you recently seen online advertisements related to 'Soccer shoes'?

- Yes
- No
- I don't remember

* Have the online ads you have seen related to 'Soccer shoes' been relevant and useful to you (e.g. the ad introduced you to a new product that you appreciate, or reminded you to purchase a product you had intended to buy)?

- Yes
- Maybe
- No

Summary of Data Collection

220 participants (82 from Pakistan, 138 from US)

For each participant, we have:

Summary of Data Collection

220 participants (82 from Pakistan, 138 from US)

For each participant, we have:

Foreground



Background

Summary of Data Collection

220 participants (82 from Pakistan, 138 from US)

For each participant, we have:

Foreground

- Survey
 1. Basic demographics
 2. General web usage
 3. Interaction with Ads
 4. Privacy practices
 5. Knowledge about APMs
 6. Relevance of interests

Background

Summary of Data Collection

220 participants (82 from Pakistan, 138 from US)

For each participant, we have:

Foreground

- Survey
 1. Basic demographics
 2. General web usage
 3. Interaction with Ads
 4. Privacy practices
 5. Knowledge about APMs
 6. Relevance of interests

Background

- Interests from 4 APMS
 1. Facebook
 2. Google
 3. BlueKai
 4. eXelate
- Browsing history (last 3 months)
- Search term history (last 3 months)

Goals of the Study

1. Who knows what and how much?
 - What inferences are drawn by each APM?
 - Does every APM infer the same information?
2. How do users perceive these interests inferred about them?

Which APM Knows More?

Which APM Knows More?

Table: Interests gathered from 220 participants

Inferred Interests				
APM	Users	Unique	Total	Avg. per User
Google	213	594	9,013	42.3
Facebook	208	25,818	108,930	523.7
BlueKai	220	3,522	92,926	422.4
eXelate	218	139	1,941	8.9

Which APM Knows More?

Table: Interests gathered from 220 participants

Inferred Interests				
APM	Users	Unique	Total	Avg. per User
Google	213	594	9,013	42.3
Facebook	208	25,818	108,930	523.7
BlueKai	220	3,522	92,926	422.4
eXelate	218	139	1,941	8.9

- Facebook gathers maximum interests, while eXelate has the least

Which APM Knows More?

Table: Interests gathered from 220 participants

Inferred Interests				
APM	Users	Unique	Total	Avg. per User
Google	213	594	9,013	42.3
Facebook	208	25,818	108,930	523.7
BlueKai	220	3,522	92,926	422.4
eXelate	218	139	1,941	8.9

- Facebook gathers maximum interests, while eXelate has the least
- Bluekai had a profile on every user

Which APM Knows More?

Table: Interests gathered from 220 participants

Inferred Interests				
APM	Users	Unique	Total	Avg. per User
Google	213	594	9,013	42.3
Facebook	208	25,818	108,930	523.7
BlueKai	220	3,522	92,926	422.4
eXelate	218	139	1,941	8.9

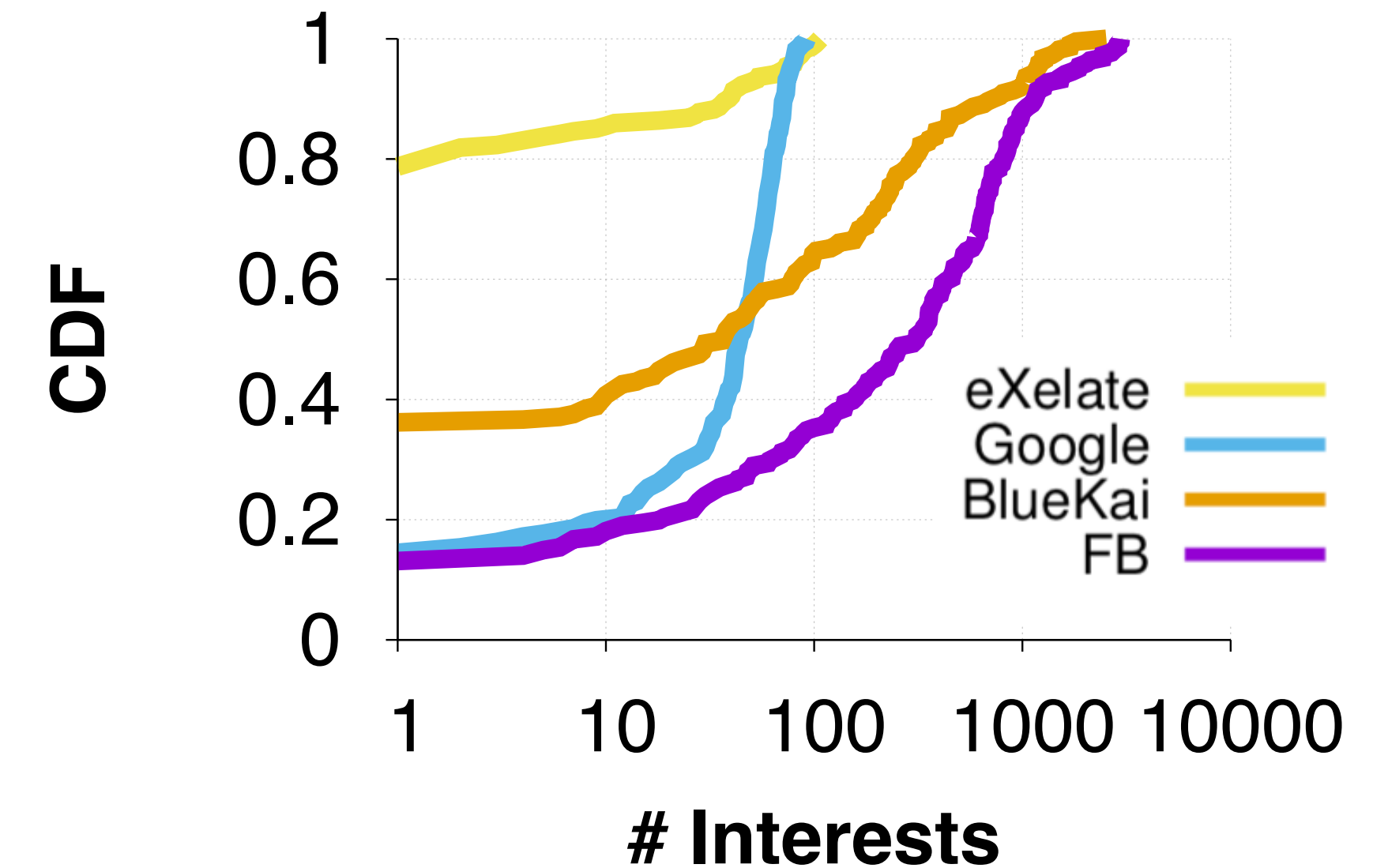


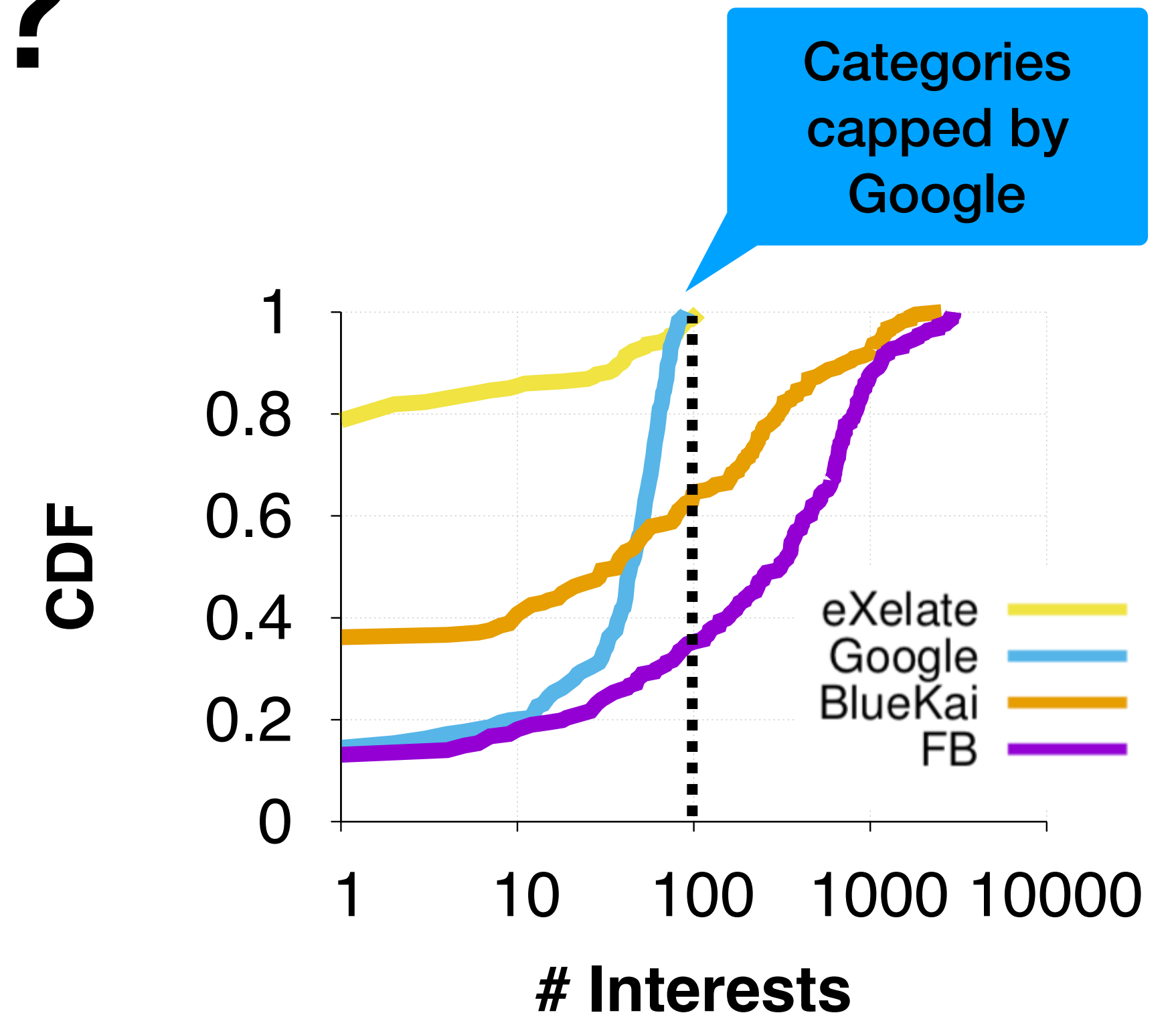
Fig: CDF of interests per user

- Facebook gathers maximum interests, while eXelate has the least
- Bluekai had a profile on every user

Which APM Knows More?

Table: Interests gathered from 220 participants

APM	Users	Inferred Interests		
		Unique	Total	Avg. per User
Google	213	594	9,013	42.3
Facebook	208	25,818	108,930	523.7
BlueKai	220	3,522	92,926	422.4
eXelate	218	139	1,941	8.9



- Facebook gathers maximum interests, while eXelate has the least
- Bluekai had a profile on every user

Fig: CDF of interests per user

Canonicalization of Interests

We cannot directly compare interests from different APMs

- **Synonyms:** Real Estate, Property
- **Granularity:** Sports, Tennis, Wimbledon

Canonicalization of Interests

We cannot directly compare interests from different APMs

- **Synonyms:** Real Estate, Property
- **Granularity:** Sports, Tennis, Wimbledon

For fair comparison, we need to map interests to a common space

Canonicalization of Interests

We cannot directly compare interests from different APMs

- **Synonyms:** Real Estate, Property
- **Granularity:** Sports, Tennis, Wimbledon

For fair comparison, we need to map interests to a common space

We used Open Directory Project (ODP)

- Manually mapped raw interest to 465 ODP categories

Canonicalization of Interests

We cannot directly compare interests from different APMs

- **Synonyms:** Real Estate, Property
- **Granularity:** Sports, Tennis, Wimbledon

For fair comparison, we need to map interests to a common space

We used Open Directory Project (ODP)

- Manually mapped raw interest to 465 ODP categories



Canonicalization of Interests

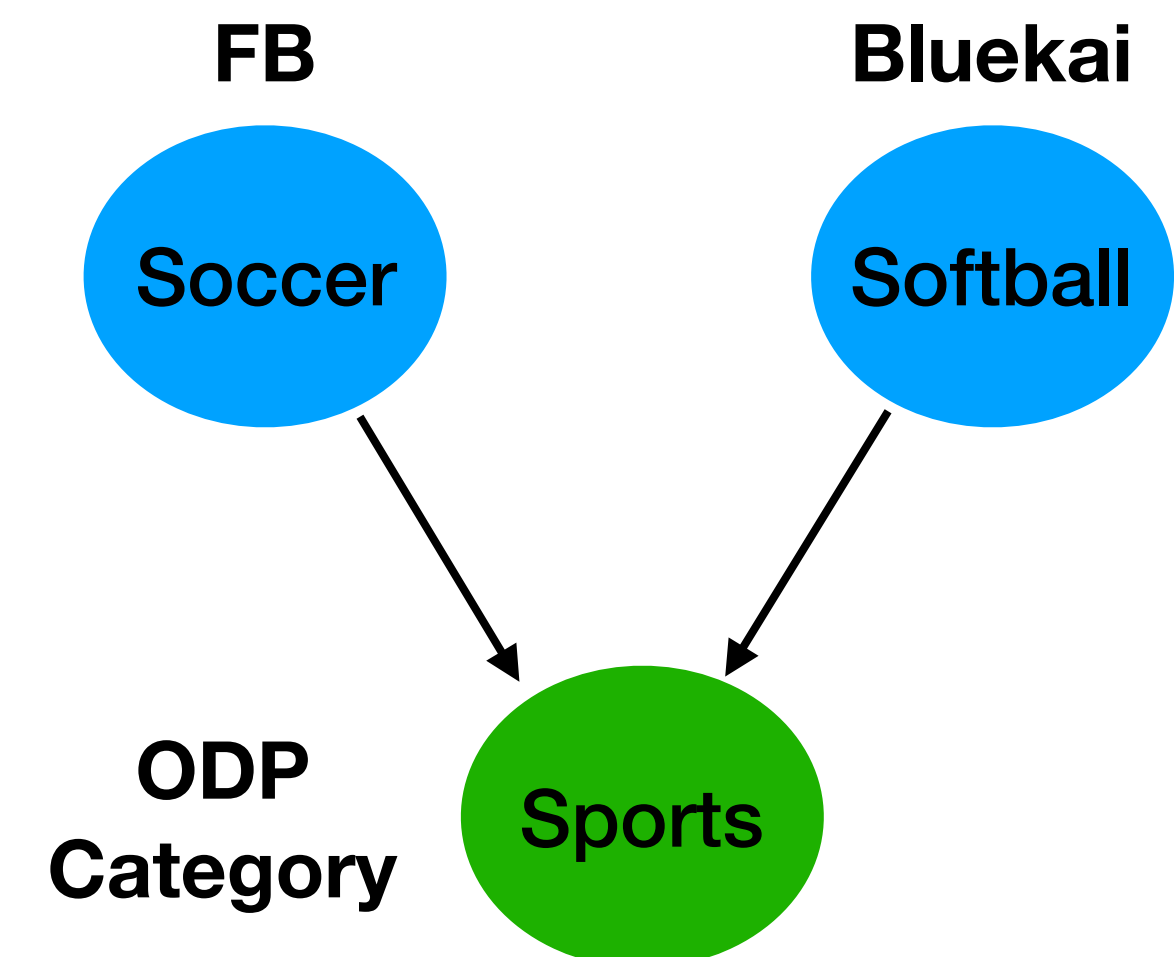
We cannot directly compare interests from different APMs

- **Synonyms:** Real Estate, Property
- **Granularity:** Sports, Tennis, Wimbledon

For fair comparison, we need to map interests to a common space

We used Open Directory Project (ODP)

- Manually mapped raw interest to 465 ODP categories



Inferred Interests After ODP Mapping

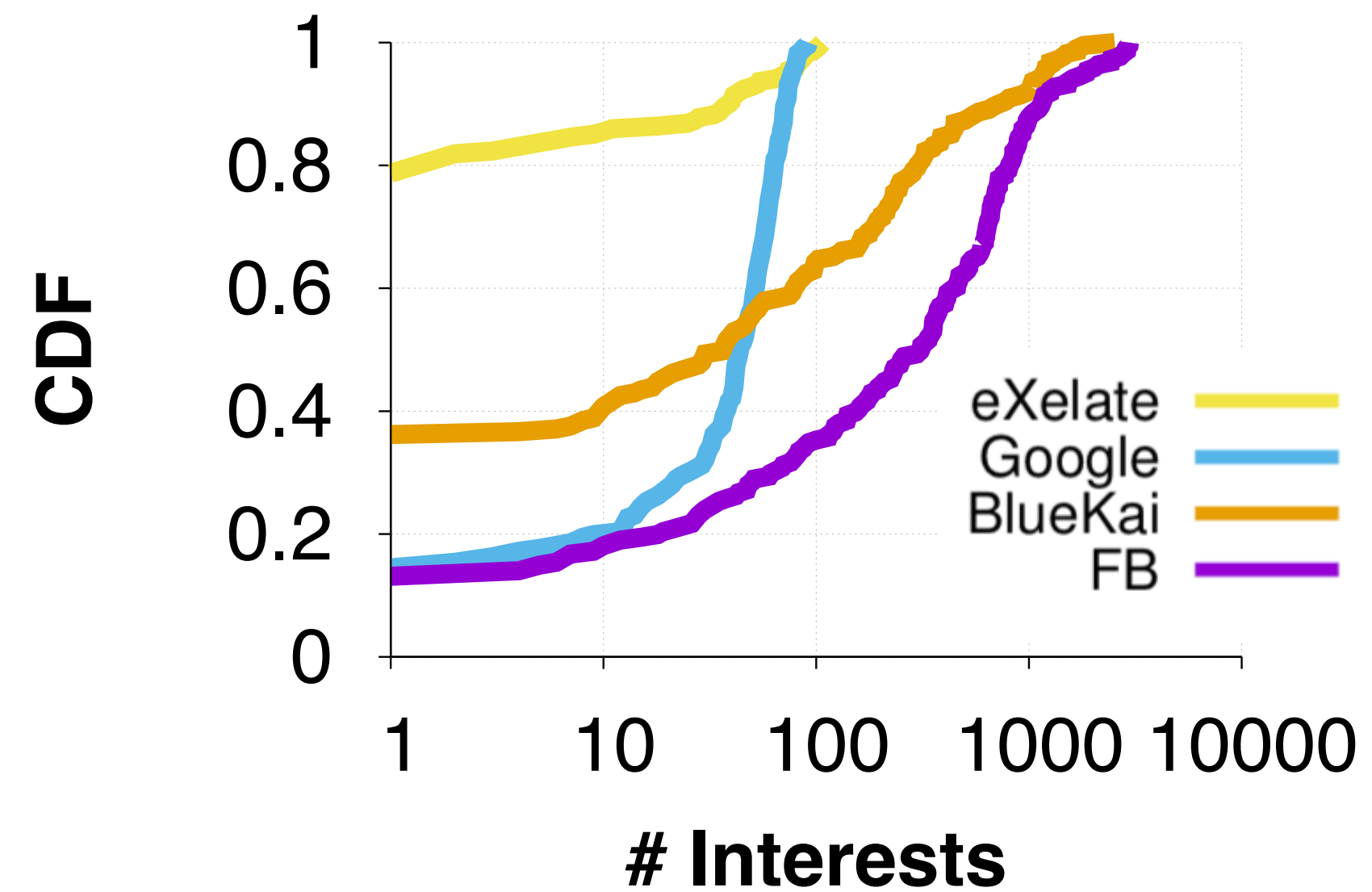


Fig: CDF of **raw** interests per user

Inferred Interests After ODP Mapping

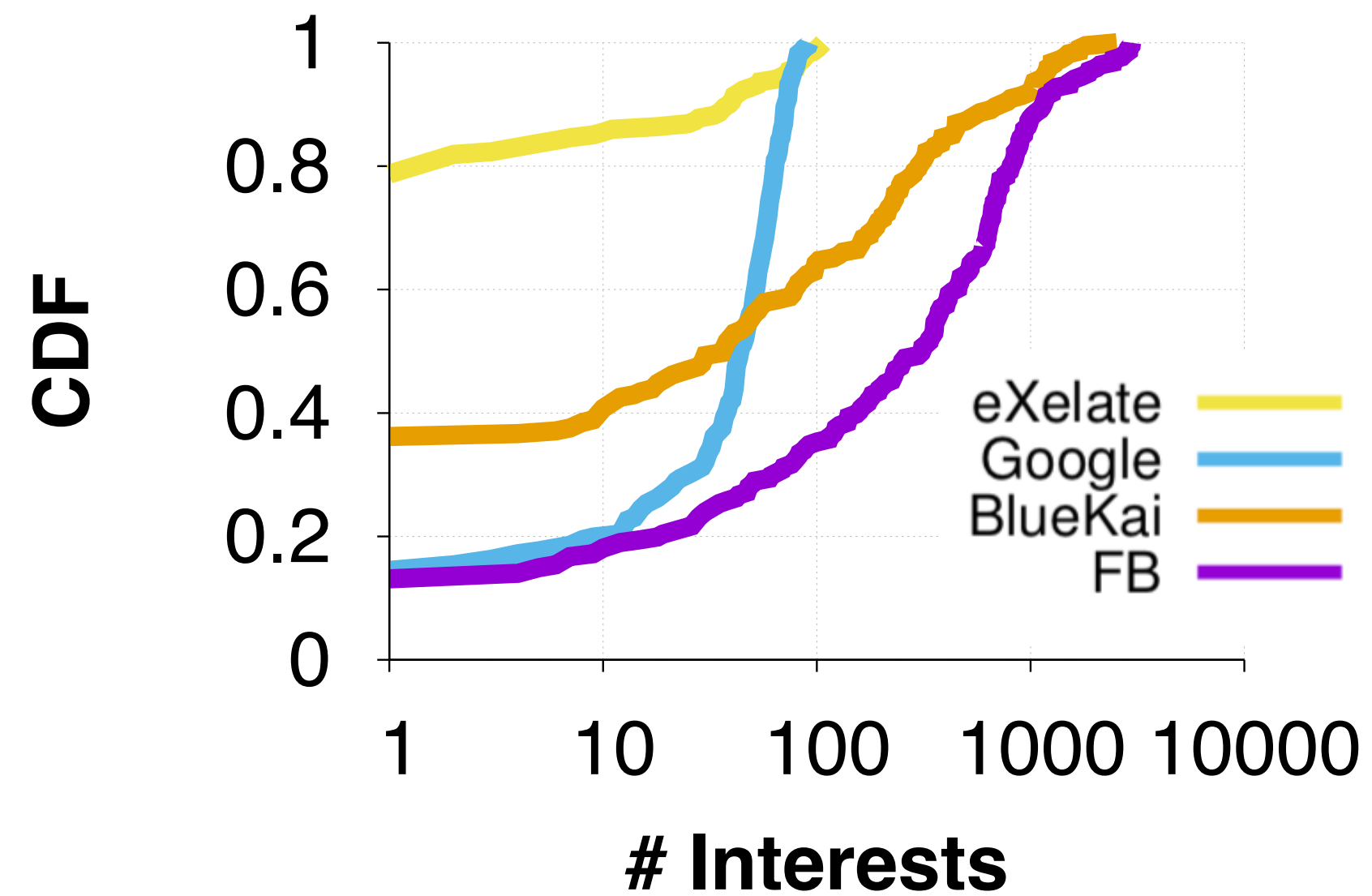


Fig: CDF of **raw** interests per user

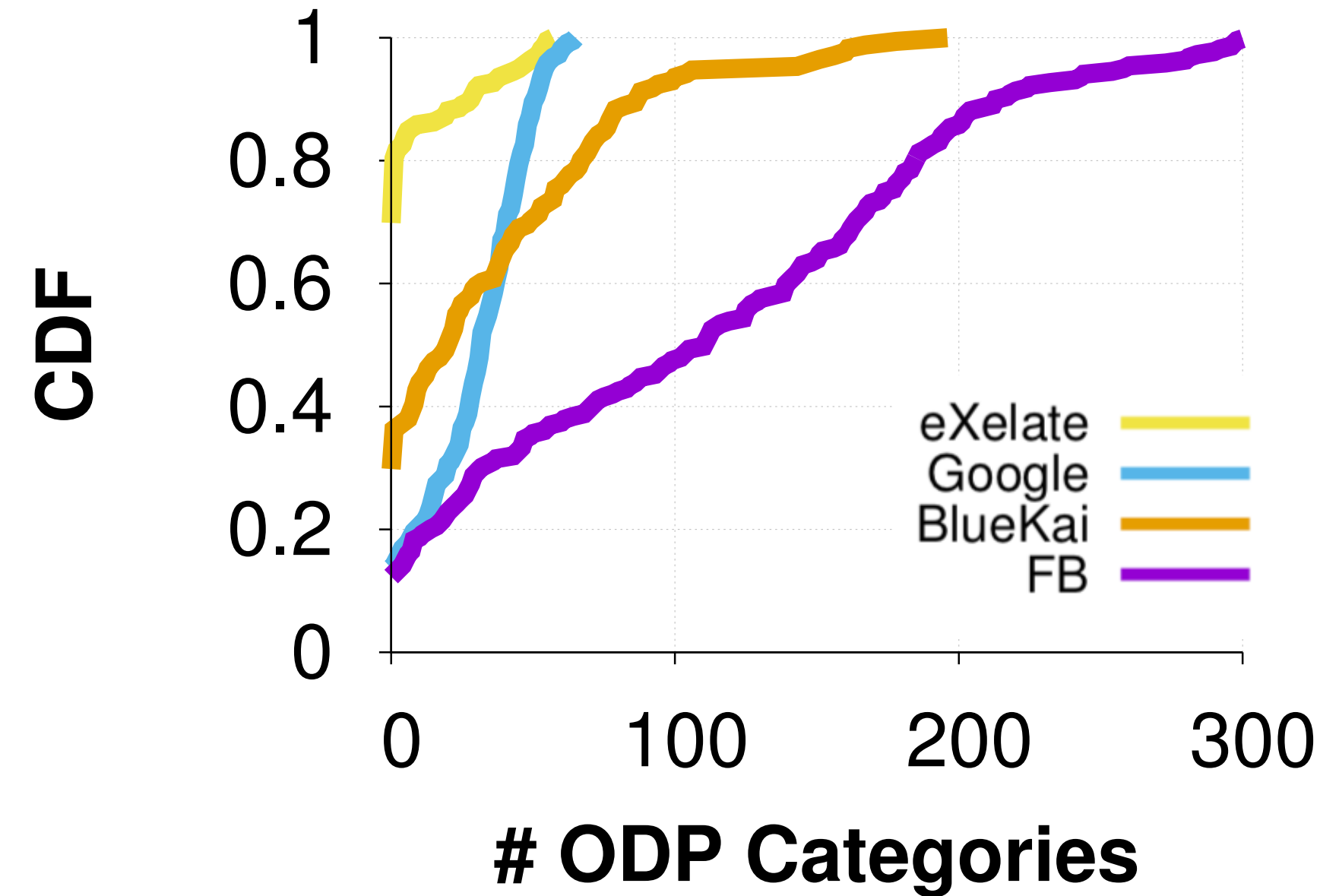


Fig: CDF of **ODP** categories per user

Do APMs Infer Similar Interests?

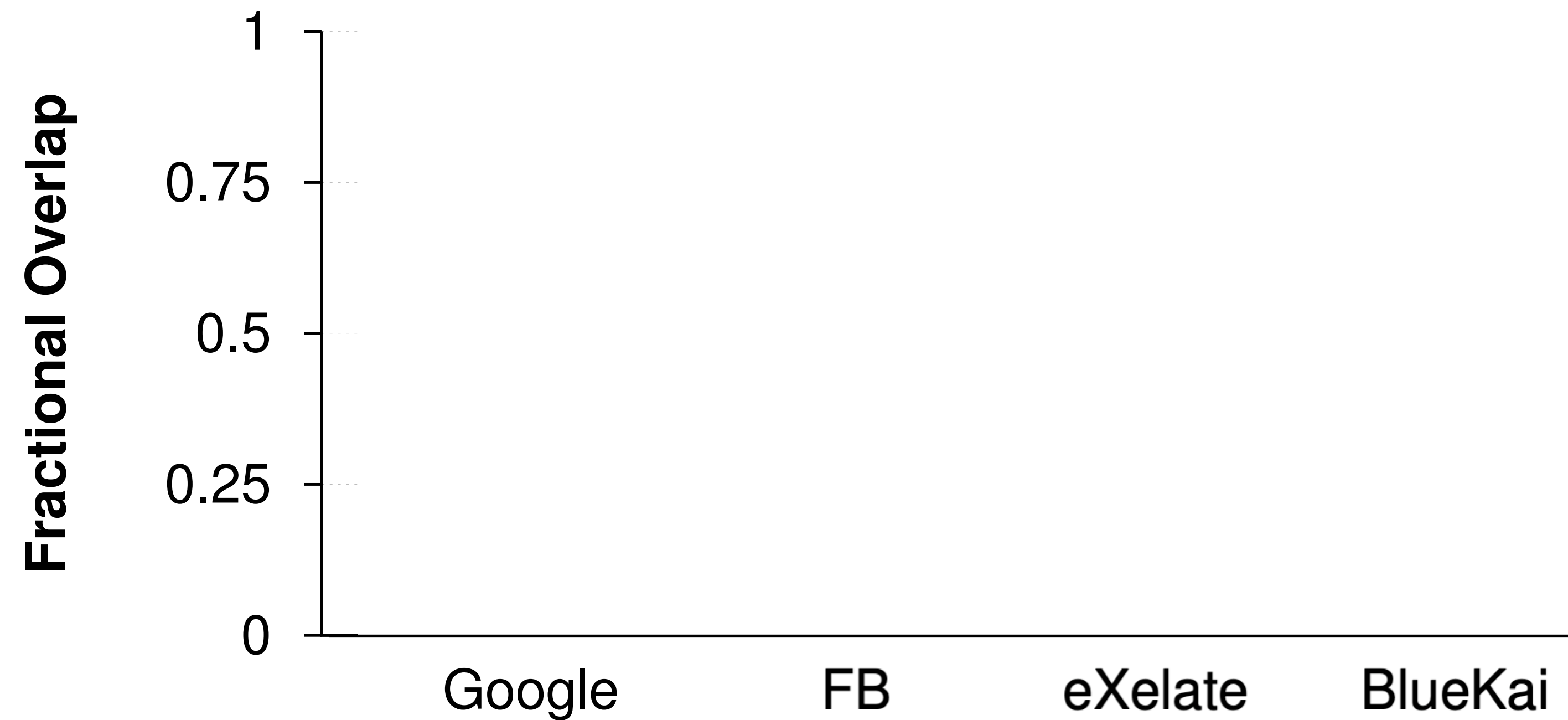


Fig: Per Participant overlap of ODP categorized interests (min, 5th, median, 95th, max)

Do APMs Infer Similar Interests?

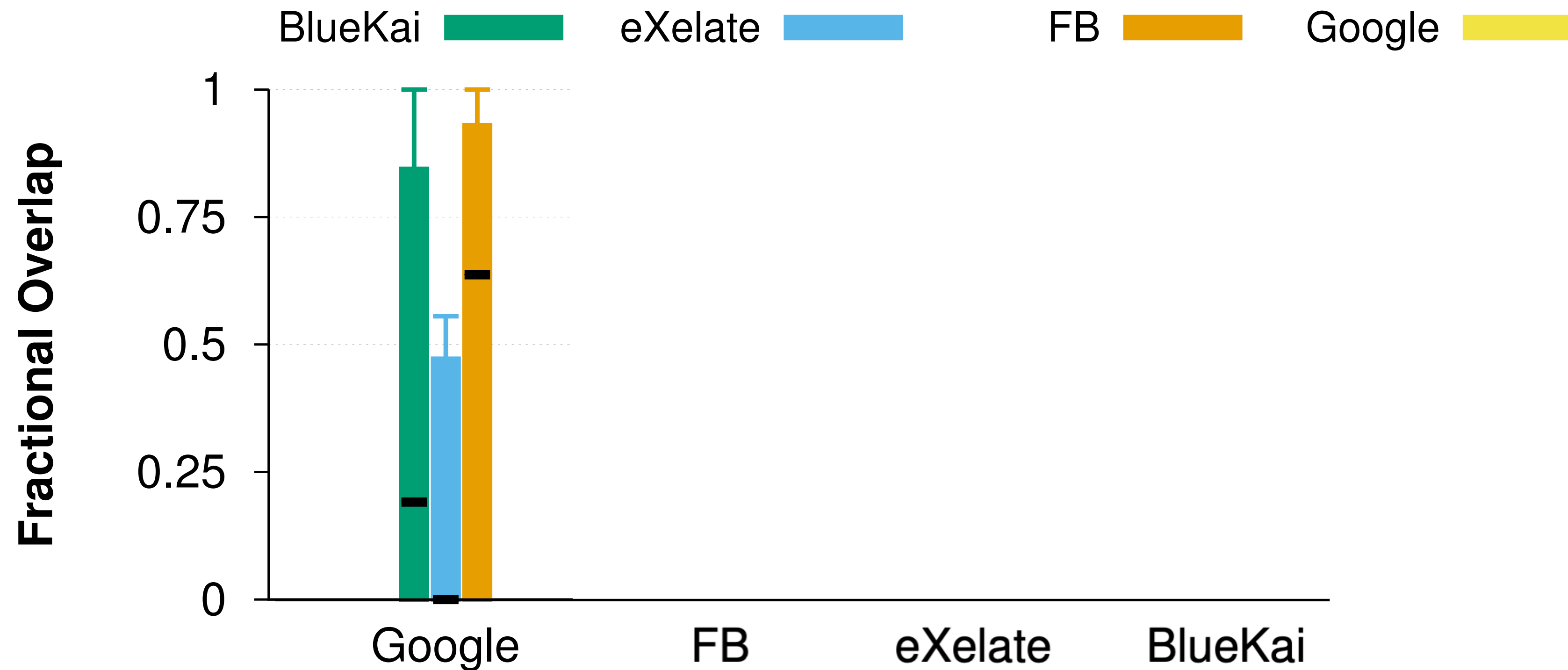


Fig: Per Participant overlap of ODP categorized interests (min, 5th, median, 95th, max)

Do APMs Infer Similar Interests?

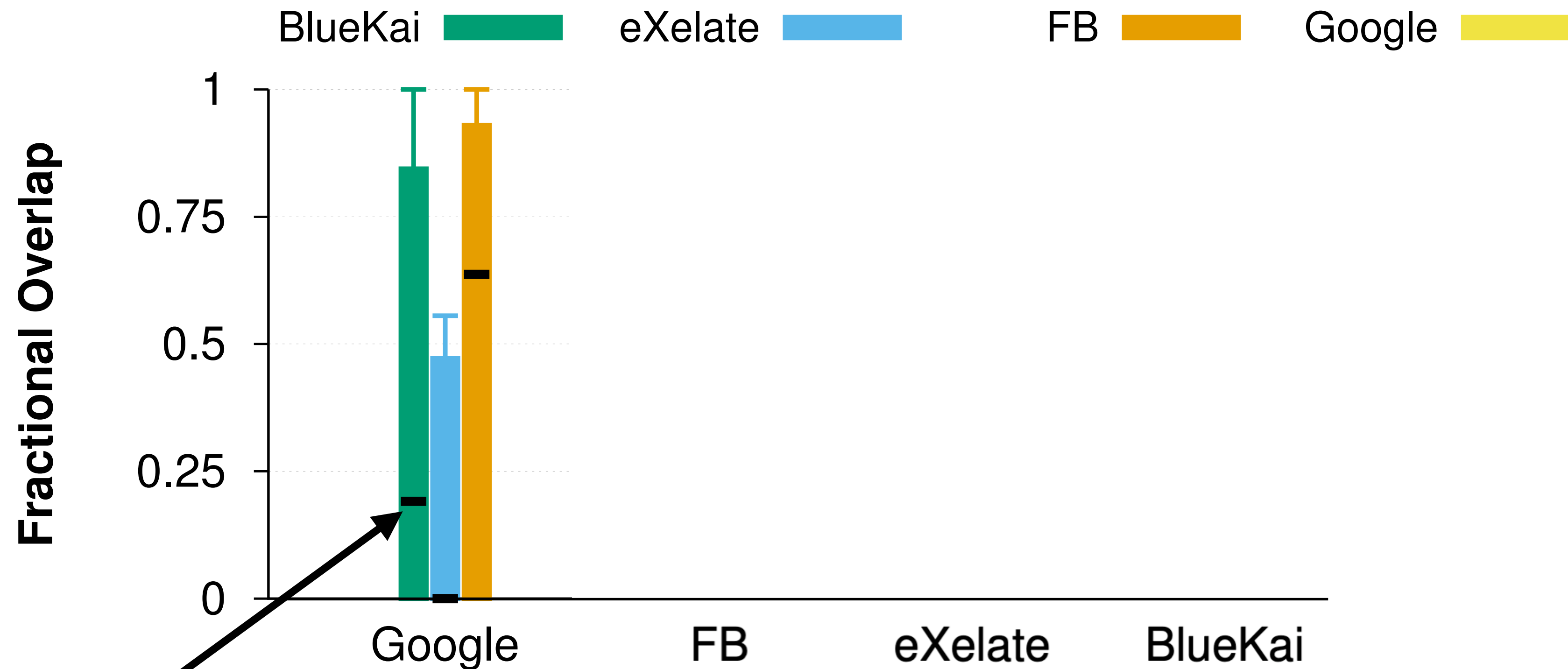


Fig: Per Participant overlap of ODP categorized interests (min, 5th, median, 95th, max)

Median Google user's interest profile has 20% overlap with BlueKai

Do APMs Infer Similar Interests?

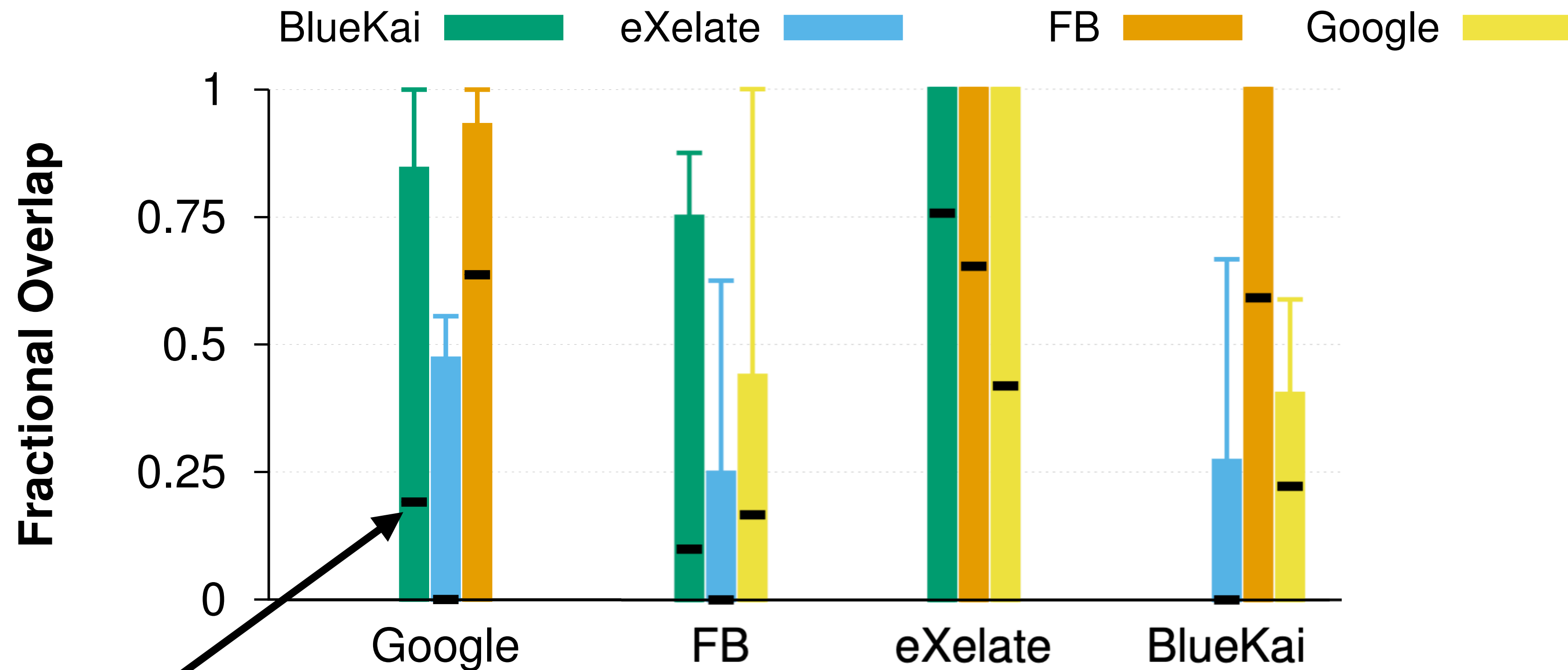


Fig: Per Participant overlap of ODP categorized interests (min, 5th, median, 95th, max)

Median Google user's interest profile has 20% overlap with BlueKai

Key Takeaways

Different APMs have different ‘portraits’ of users

Lack of overlap across APMs

Goals of the Study

1. Who knows what and how much?
 - What inferences are drawn by each APM?
 - Does everyone infer the same information?
2. How do users perceive these interests inferred about them?
 - Do some APMs infer more relevant interests?
 - Do users find ads targeted against these interests relevant?

“Half the money I spend on advertising is wasted; the trouble is I don't know which half.”

-- John Wanamaker

Relevant Interests According to Participants

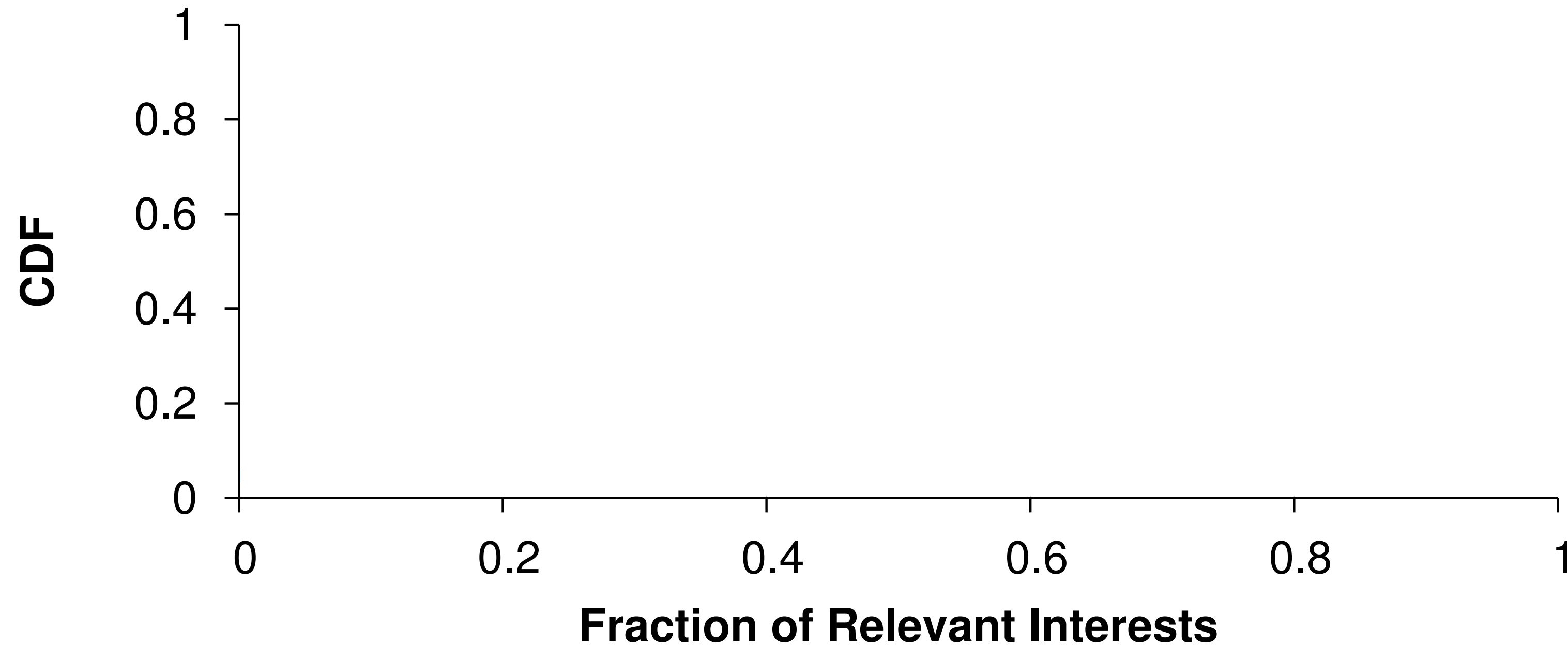


Fig: Fractions of interests rated as relevant (on a 1-5 scale) by participants

Relevant Interests According to Participants

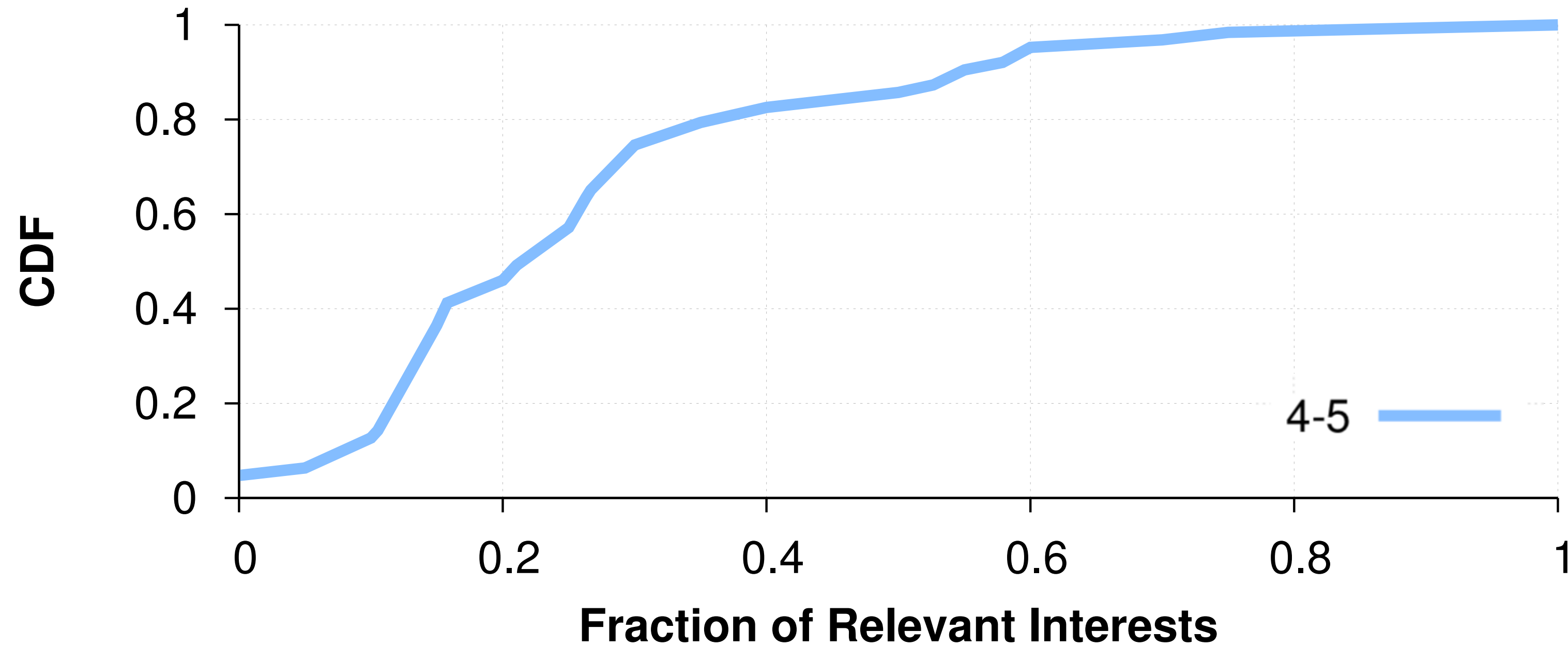


Fig: Fractions of interests rated as relevant (on a 1-5 scale) by participants

Relevant Interests According to Participants

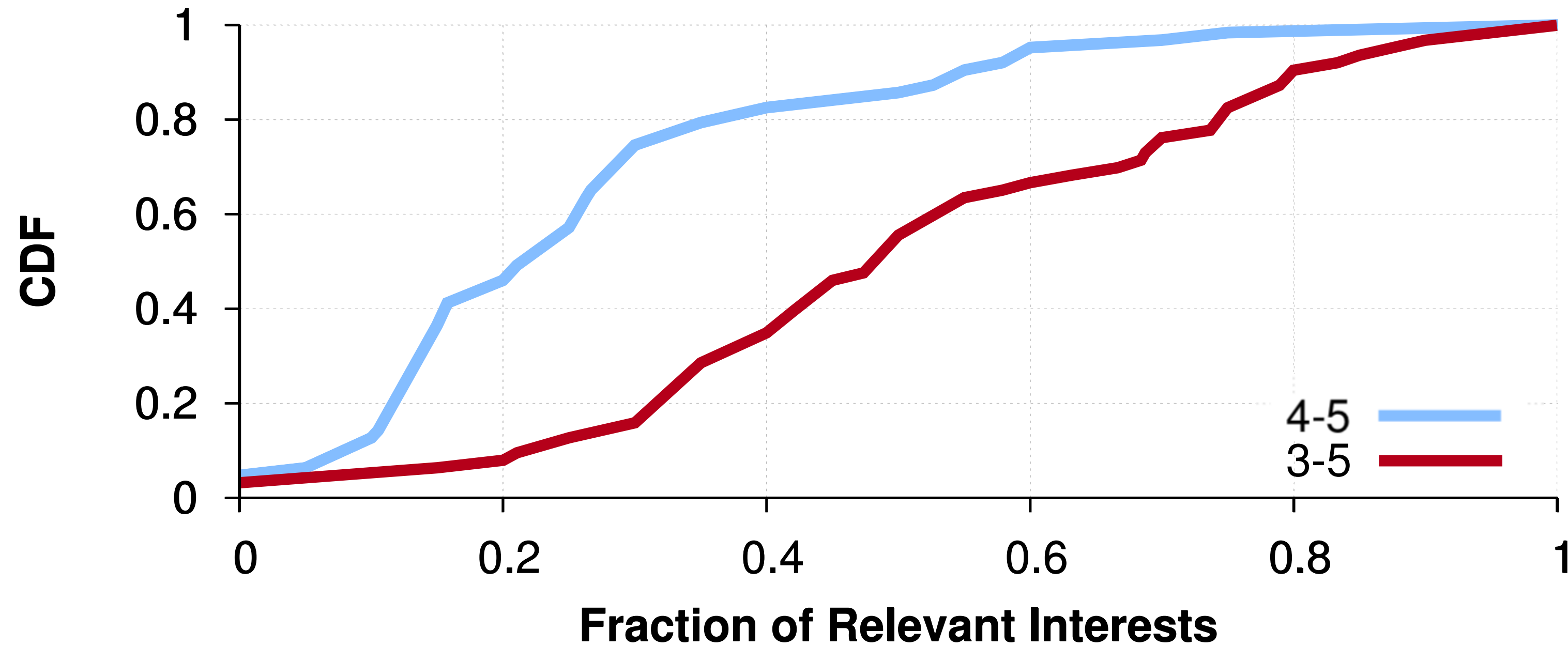


Fig: Fractions of interests rated as relevant (on a 1-5 scale) by participants

Relevant Interests According to Participants

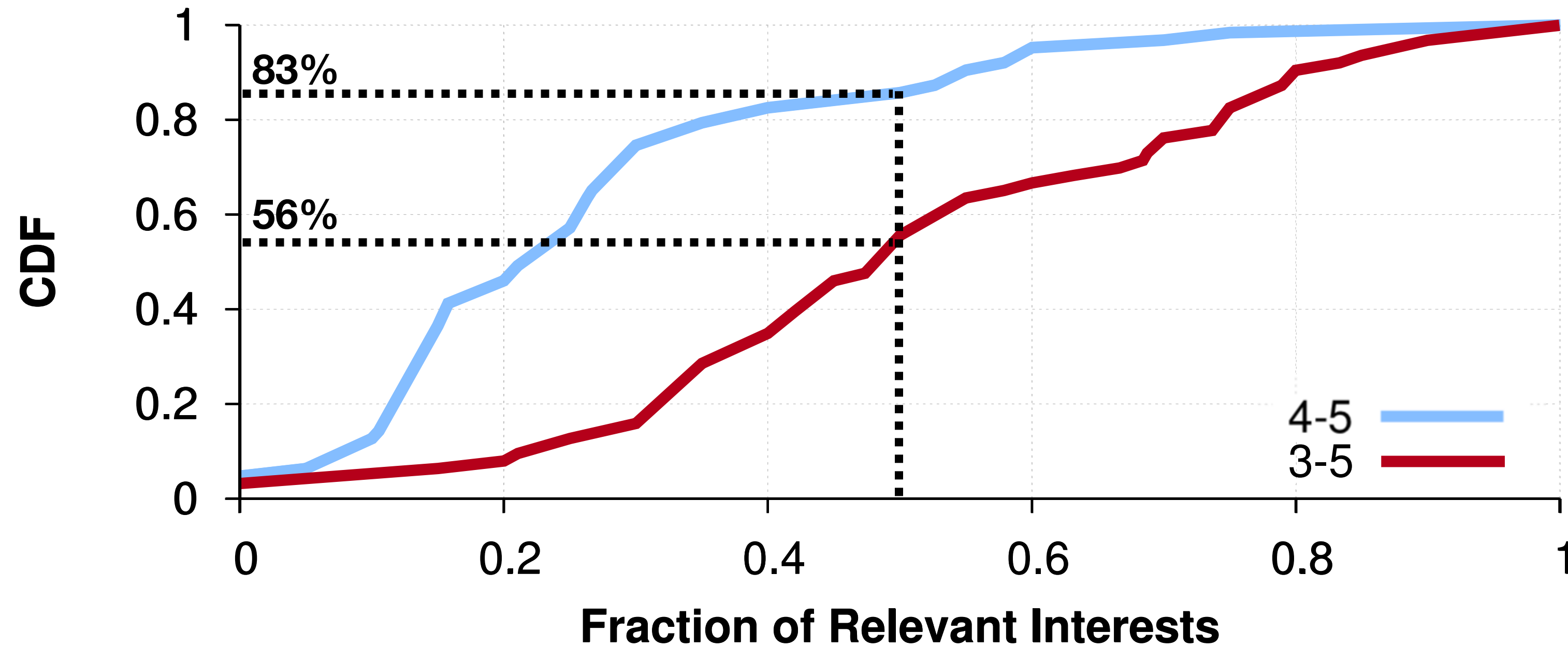


Fig: Fractions of interests rated as relevant (on a 1-5 scale) by participants

Participants' Ratings of Interests

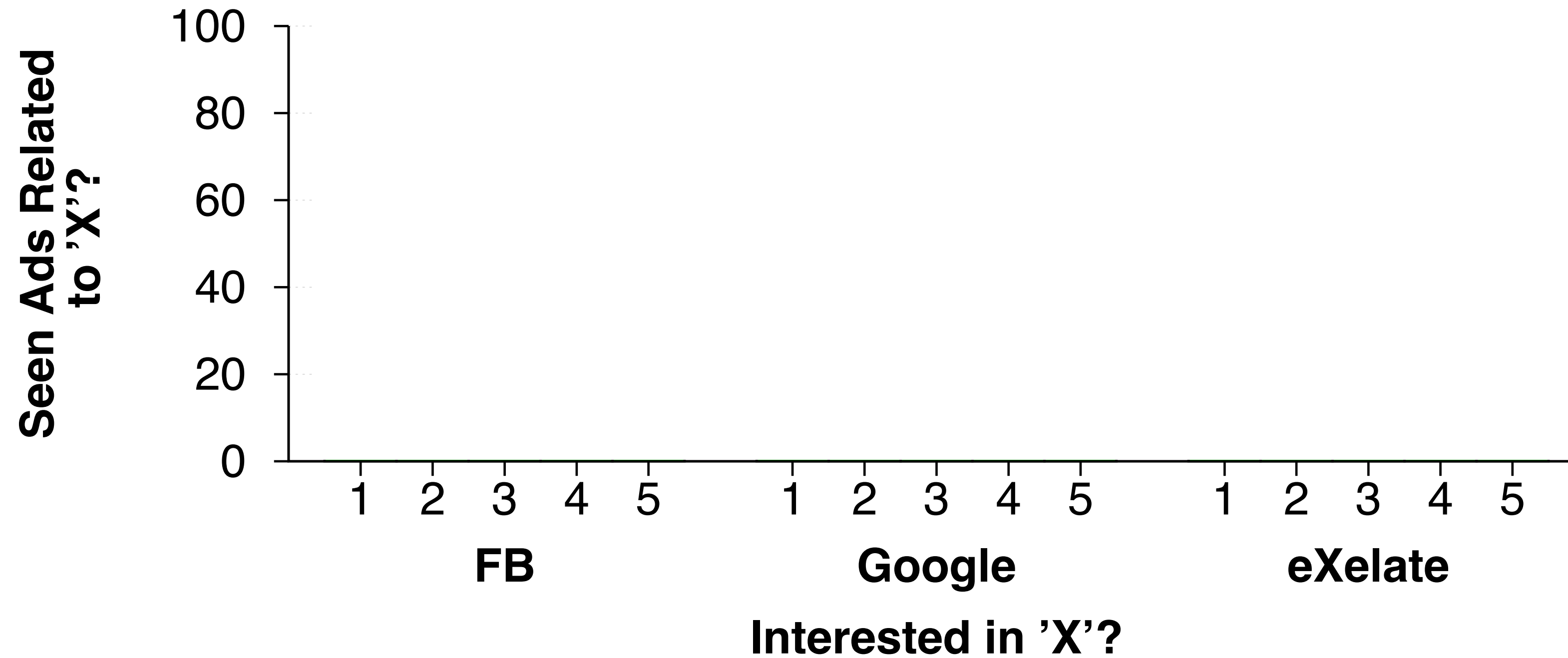


Fig: Interest Relevance vs. Seeing Ads

Participants' Ratings of Interests

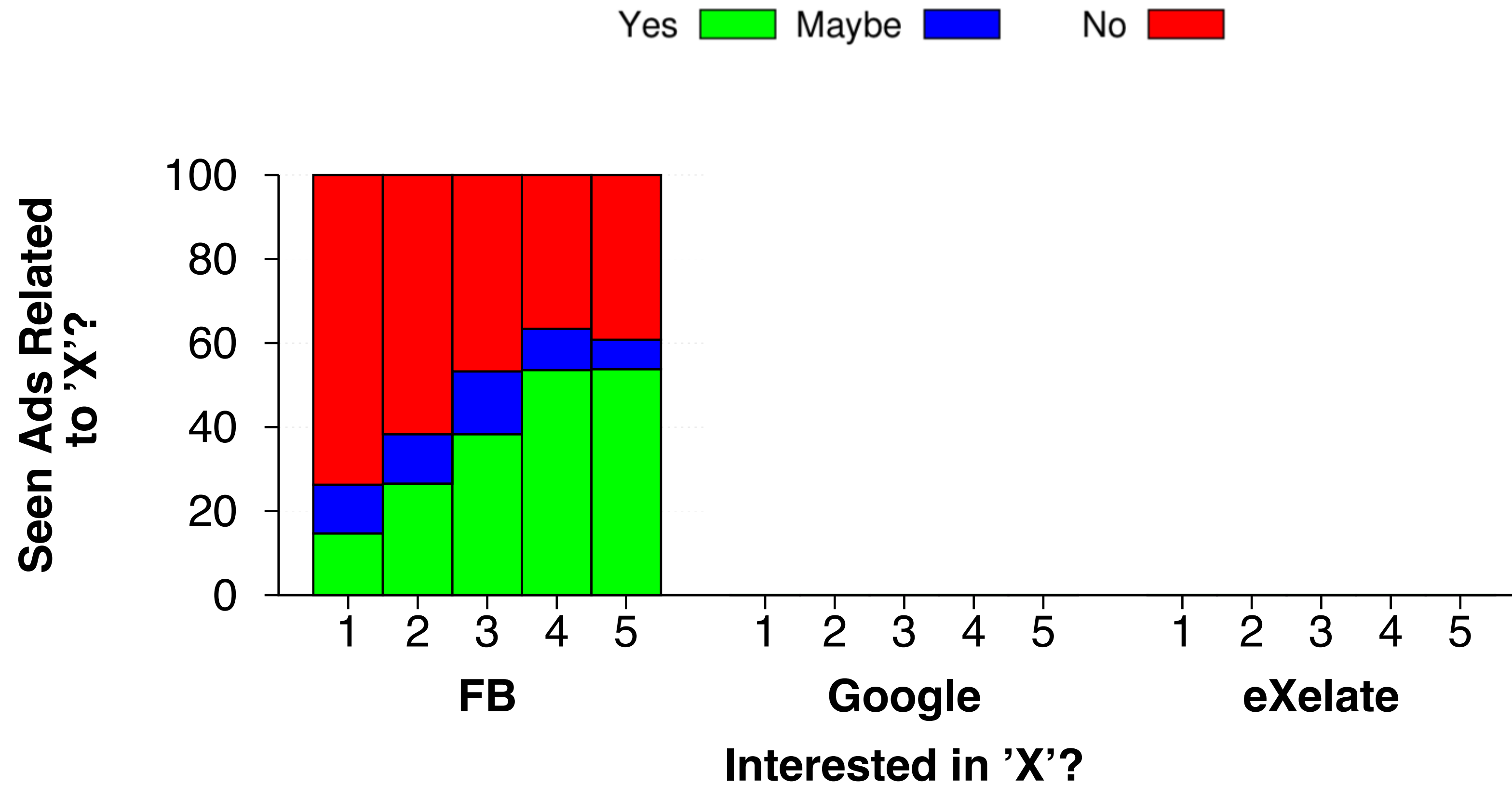


Fig: Interest Relevance vs. Seeing Ads

Participants' Ratings of Interests

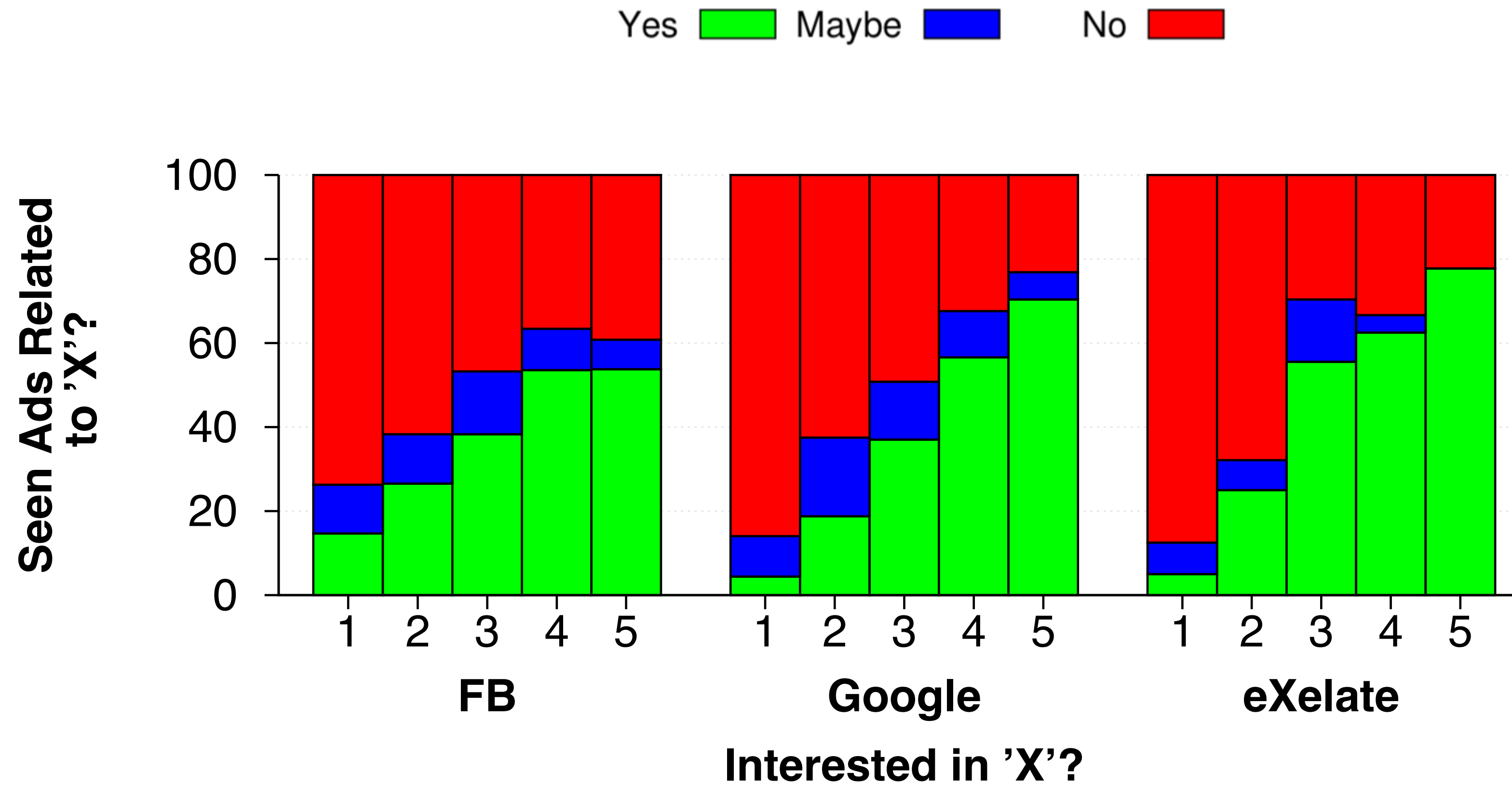


Fig: Interest Relevance vs. Seeing Ads

Participants' Ratings of Interests

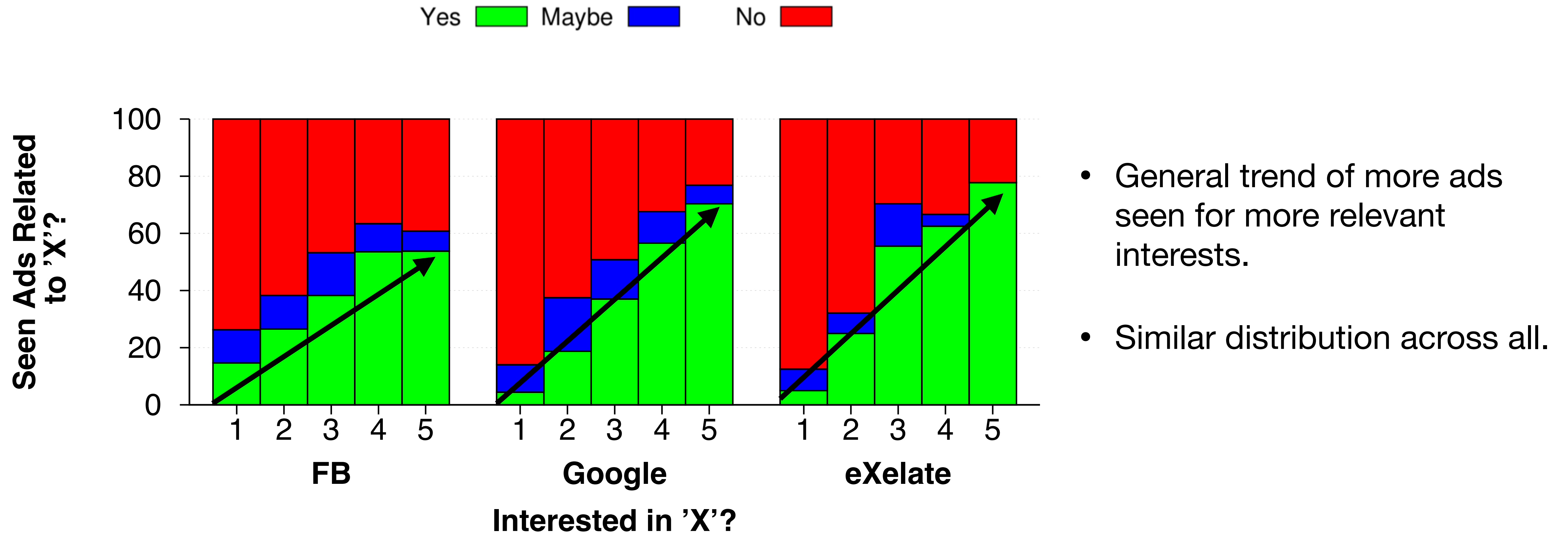


Fig: Interest Relevance vs. Seeing Ads

Participants' Ratings of Interests

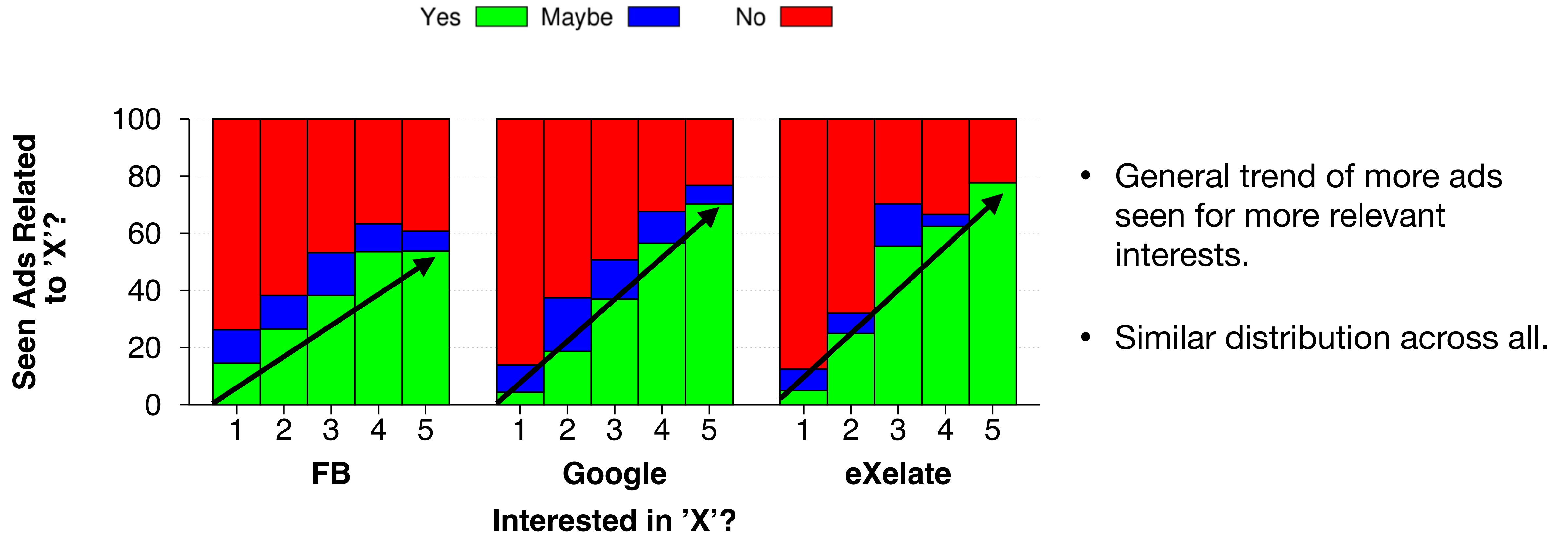


Fig: Interest Relevance vs. Seeing Ads

Majority of Interests Marked Irrelevant

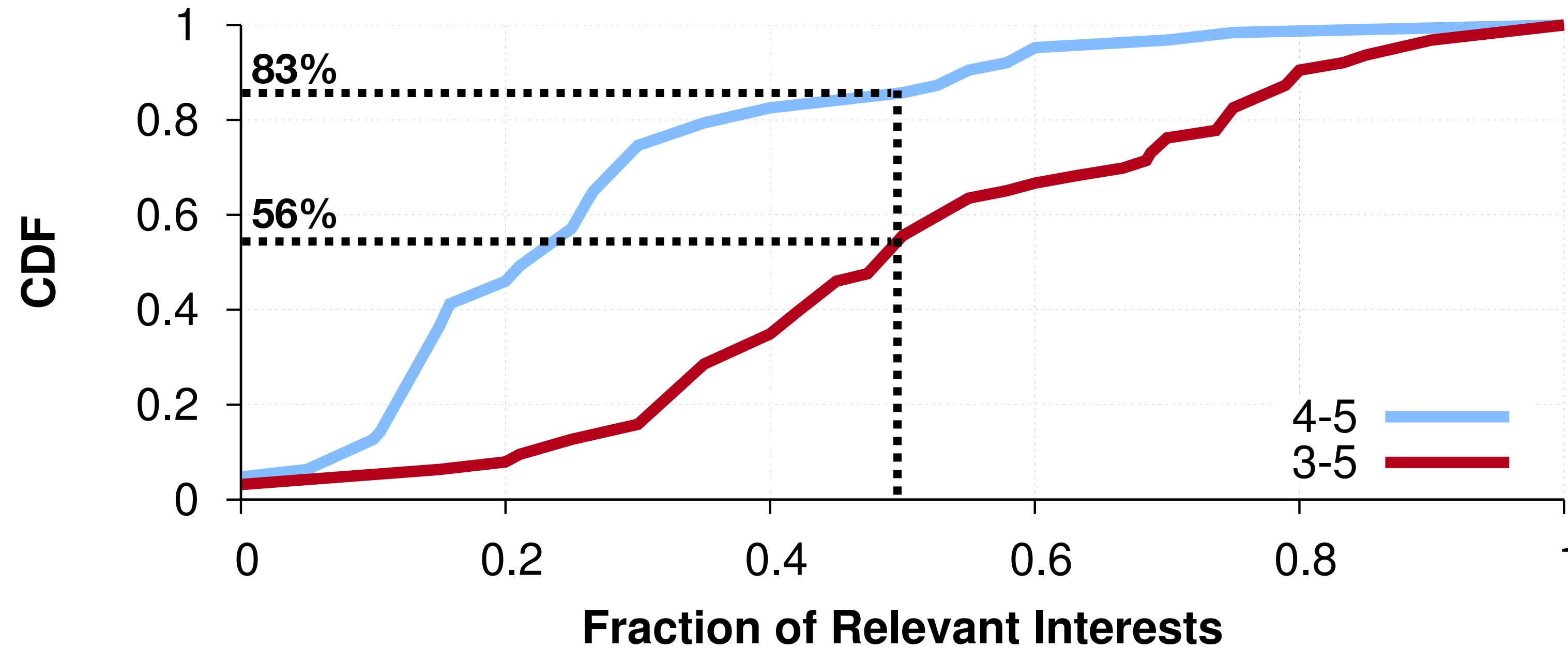


Fig: Fractions of interests rated as relevant (on a 1-5 scale) by participants

Majority of Interests Marked Irrelevant

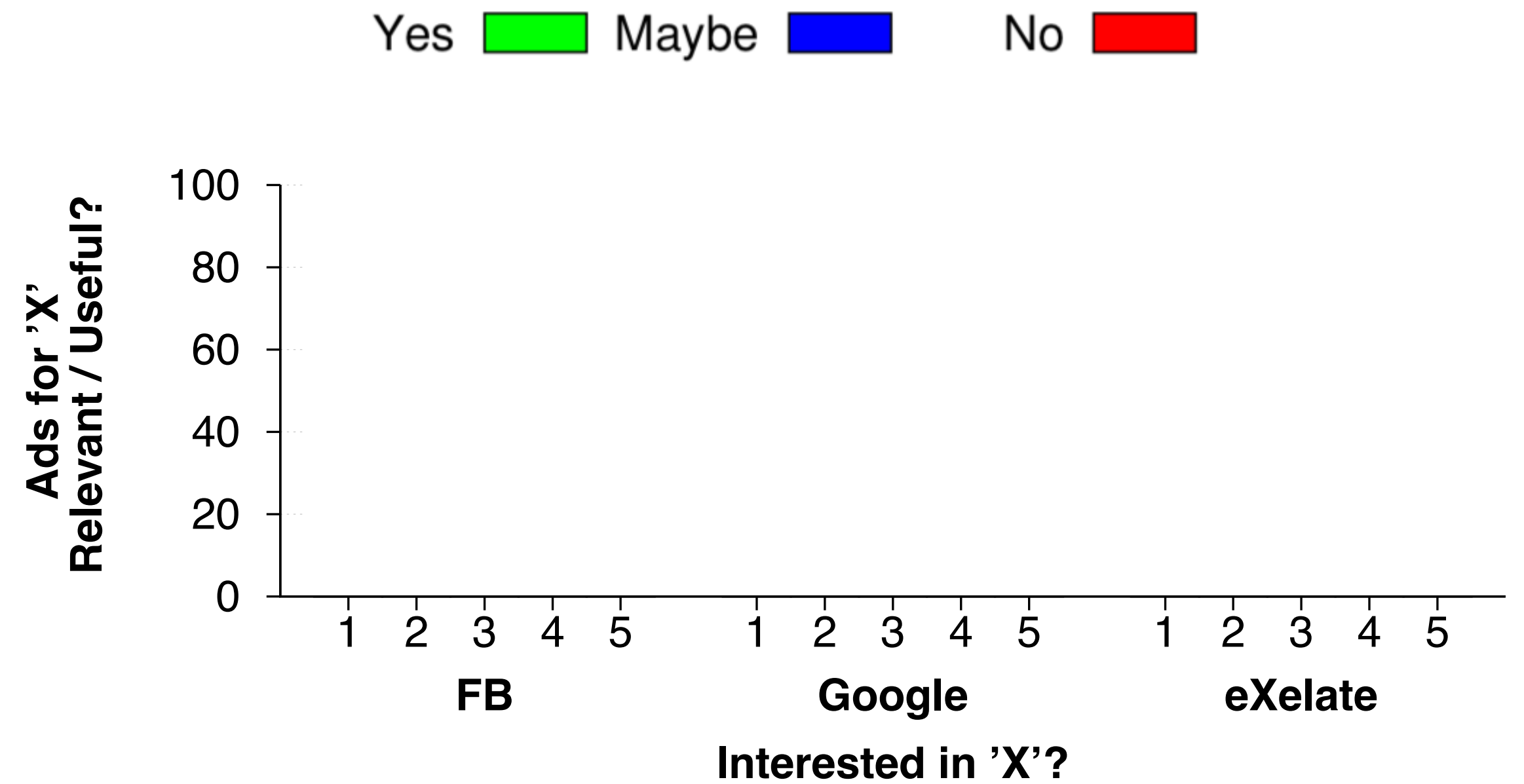
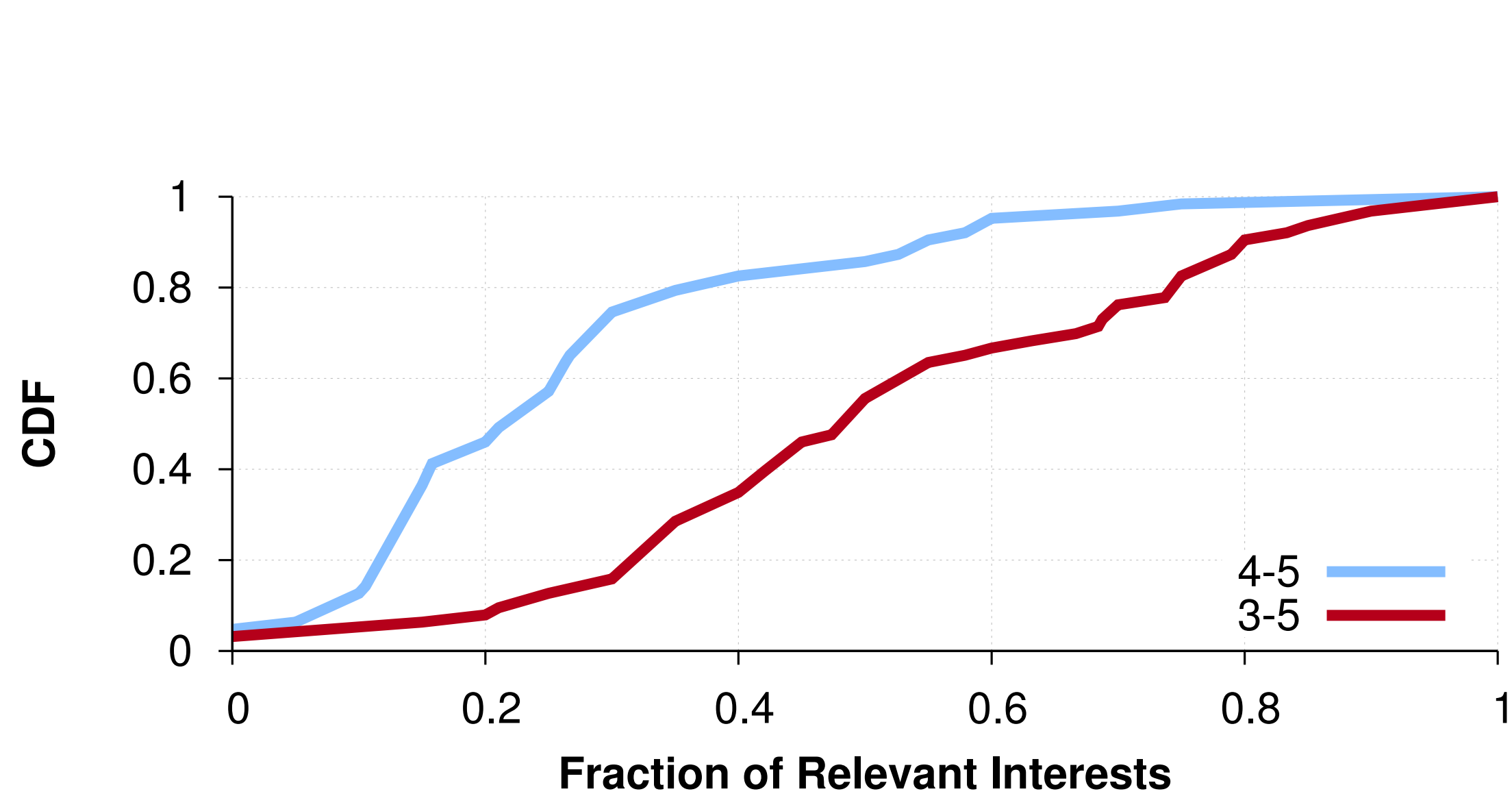


Fig: Interest Relevance vs. Seeing Relevant Ads

Majority of Interests Marked Irrelevant

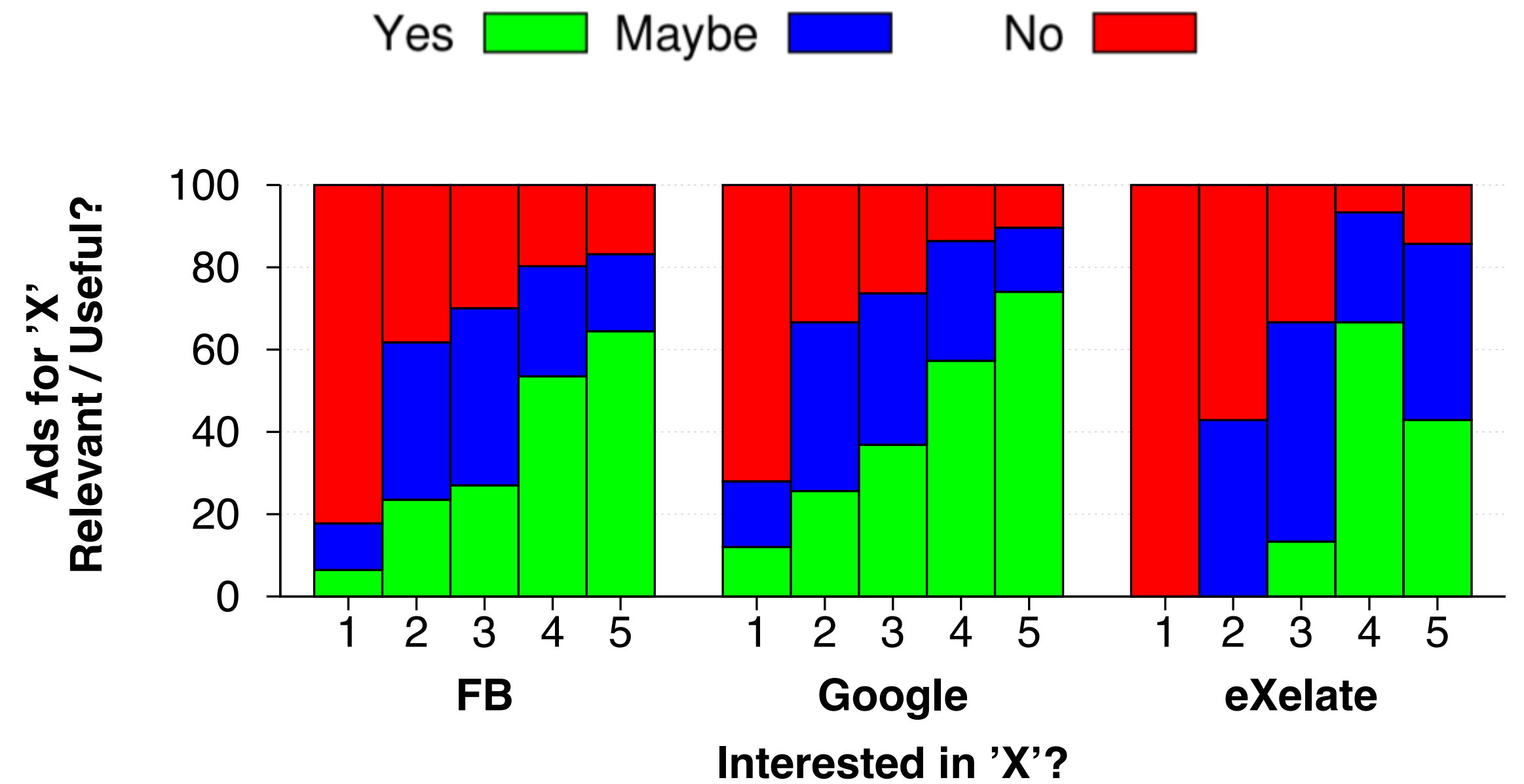
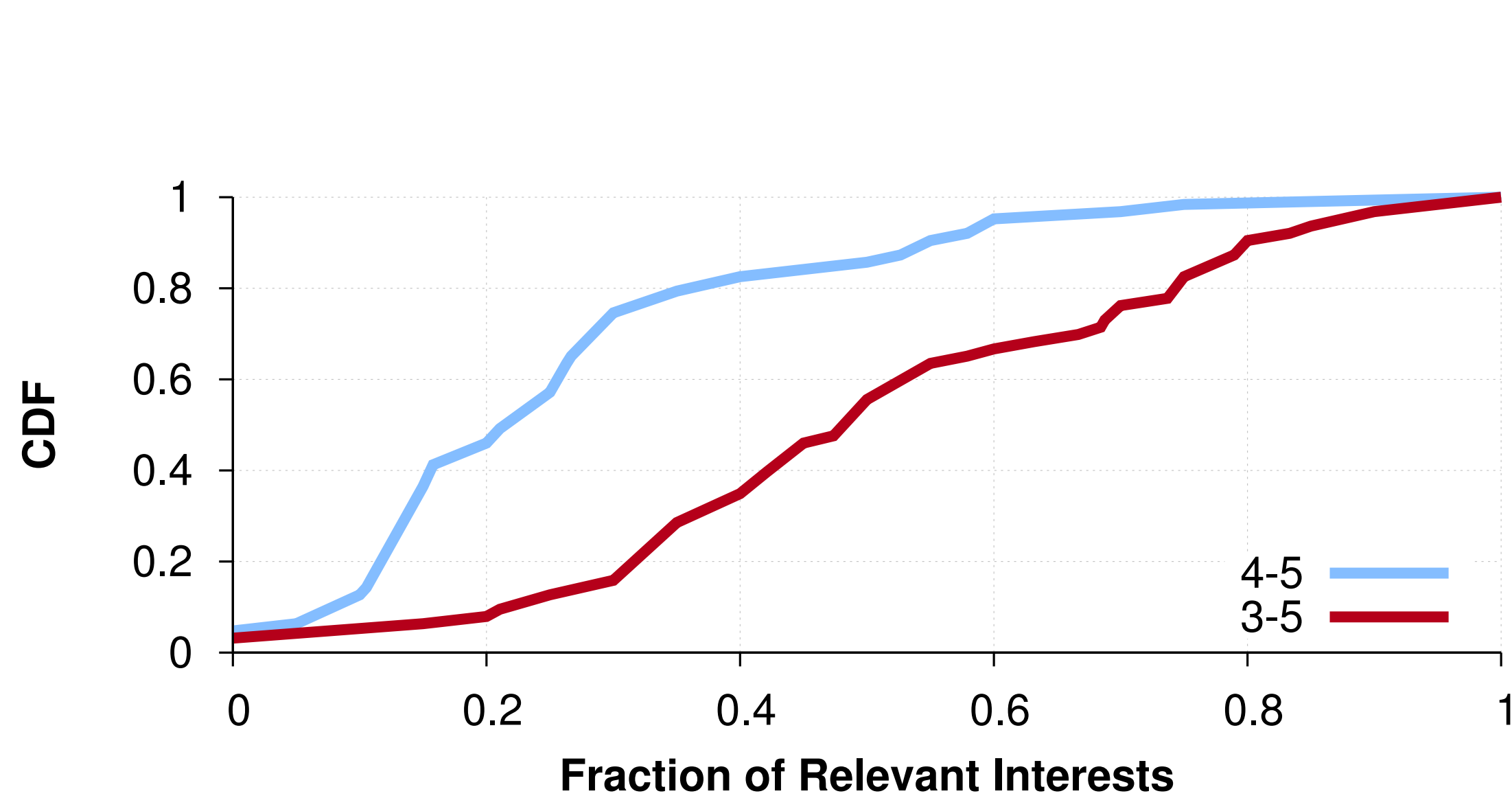


Fig: Interest Relevance vs. Seeing Relevant Ads

Majority of Interests Marked Irrelevant

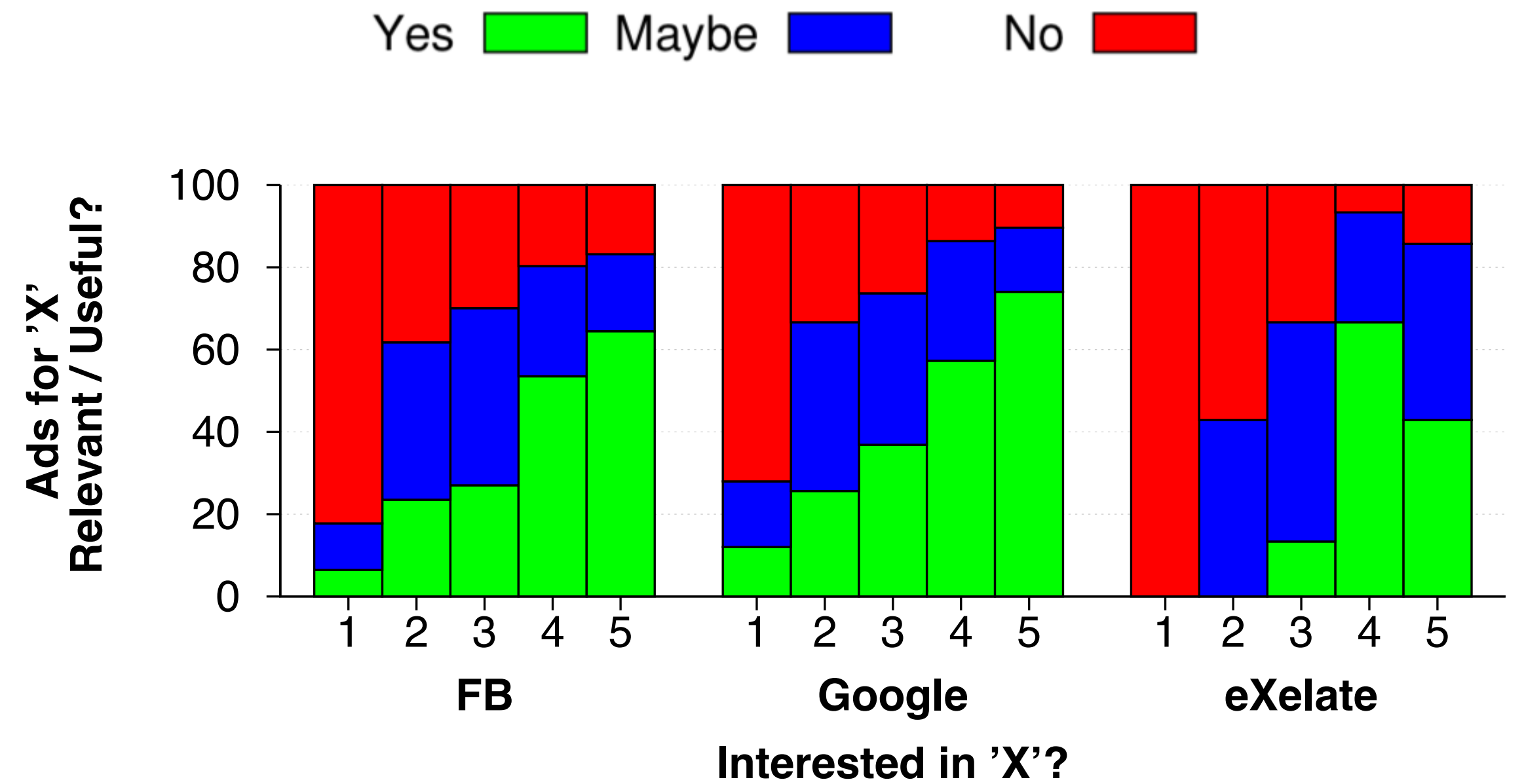
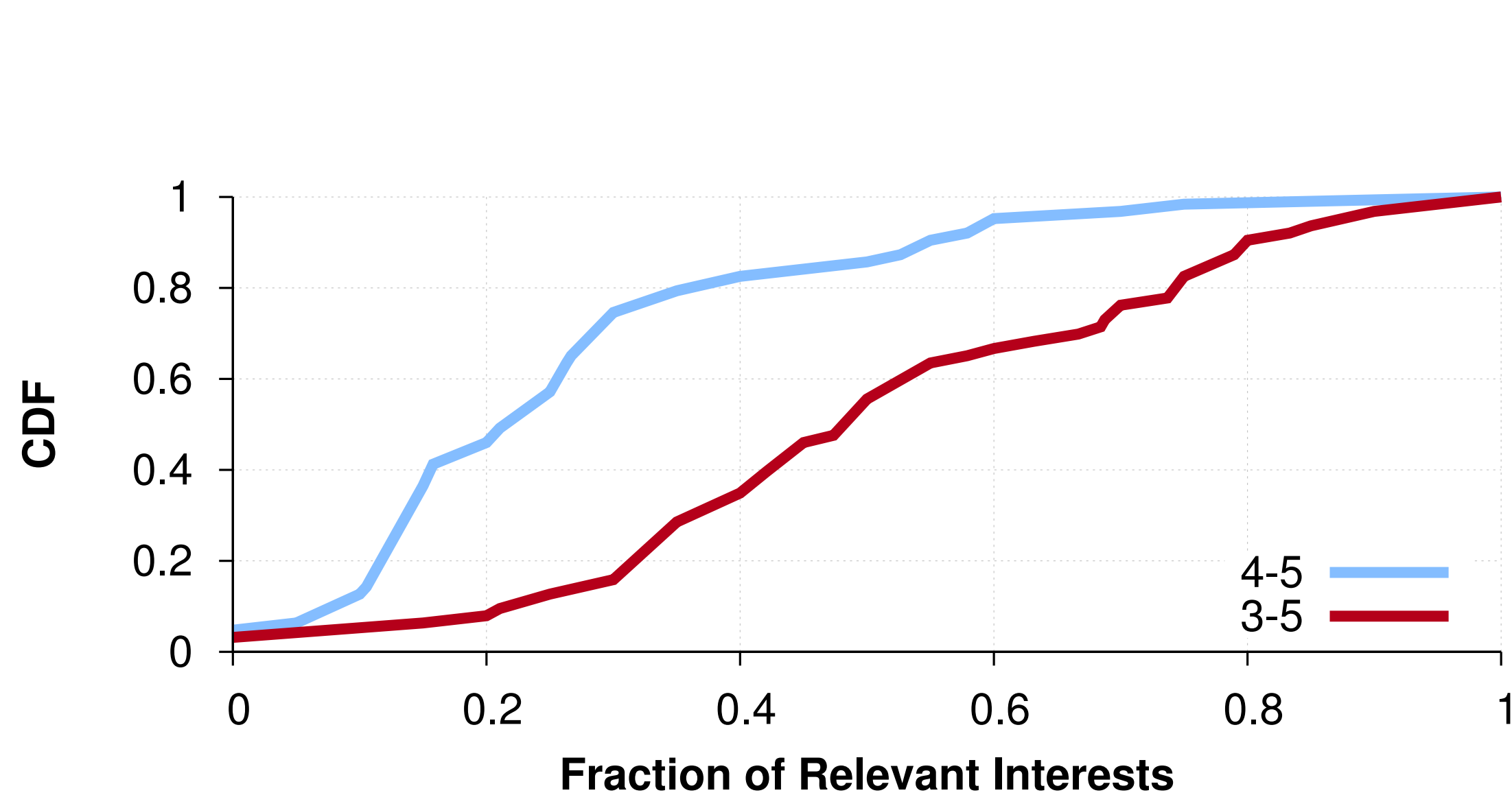


Fig: Interest Relevance vs. Seeing Relevant Ads

Users marked ads targeted to low relevant interests less useful

Key Takeaways

Majority of the interests marked not relevant

**Ads targeted to low relevance interests marked
not useful**

Limitations & Challenges

1. Participant sample is not representative of all web users
2. Single snapshot of APMs.
 - A better way would be to conduct a longitudinal study.
3. Users can have biases in recalling relevant ads.

Summary

- First large-scale study of interest profiles from four APMs
- Different APMs have different 'portraits' of the user.
- Participants rated only $< 30\%$ interests as strongly relevant.

Q: Are the marginal utility gains from targeted ads justified at the cost of privacy?

More Results in the Paper ...

1. Origin of Interests

- What fraction of the interests could be explained by historical data?
- A majority of interests could not be explained by recent browsing history

2. Affect of privacy-conscious behaviors on interest profiles

- No significant correlations

Quantity vs. Quality: Evaluating User Interest Profiles Using Ad Preference Managers

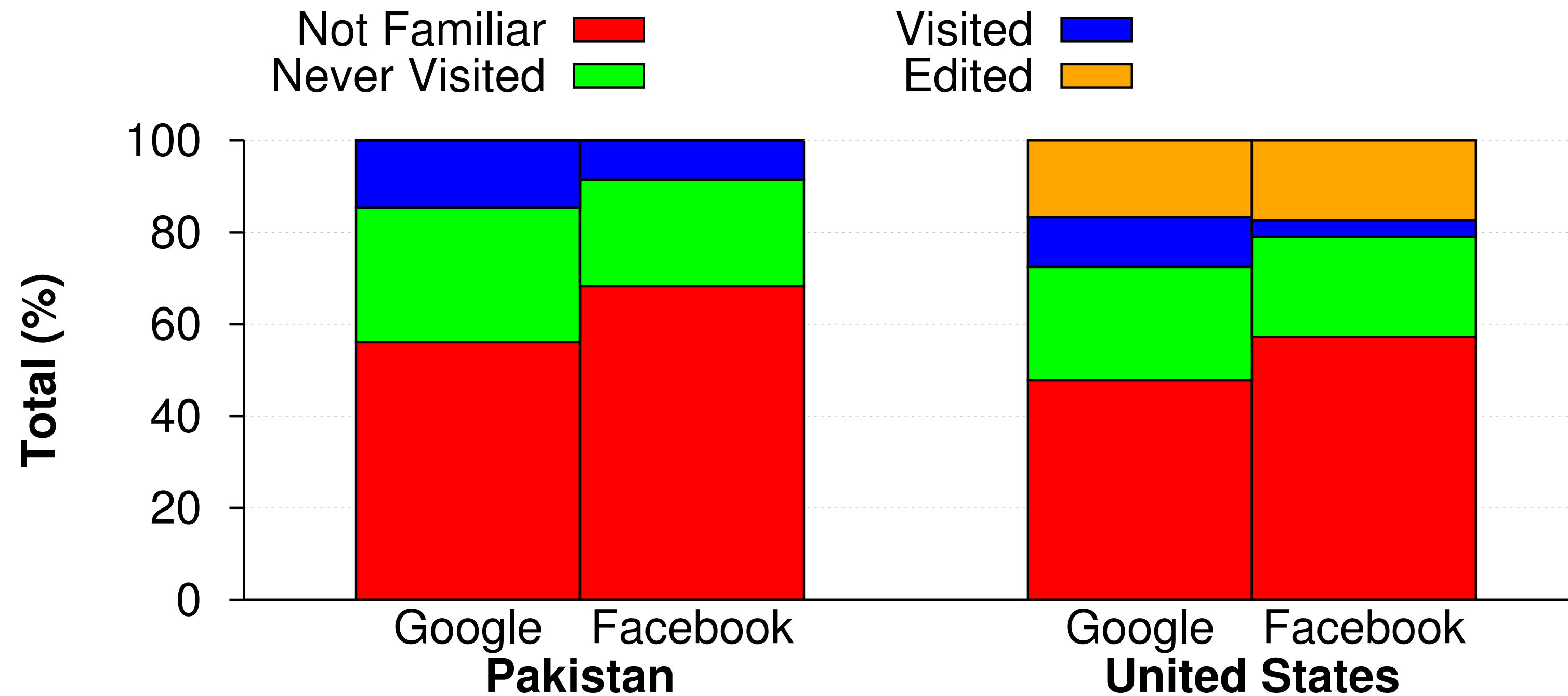
Questions?
ahmad@ccs.neu.edu

Backup Slides

Participants Dropping Out

- Overall 9 participants refused to take the survey
 - 3 provided feedback.
 - 1 did not have time and 2 had privacy reservations

Knowledge of APMs



Goals of the Study

1. Who knows what and how much?
 - What inferences are drawn by the APMs?
 - Does everyone infer the same information?
2. How do users perceive these interests inferred about them?
 - Do some APMs draw better inferences?
3. How are the inferences drawn?

How Are The Inferences Drawn?

How Are The Inferences Drawn?

Browsing History

Search History

How Are The Inferences Drawn?

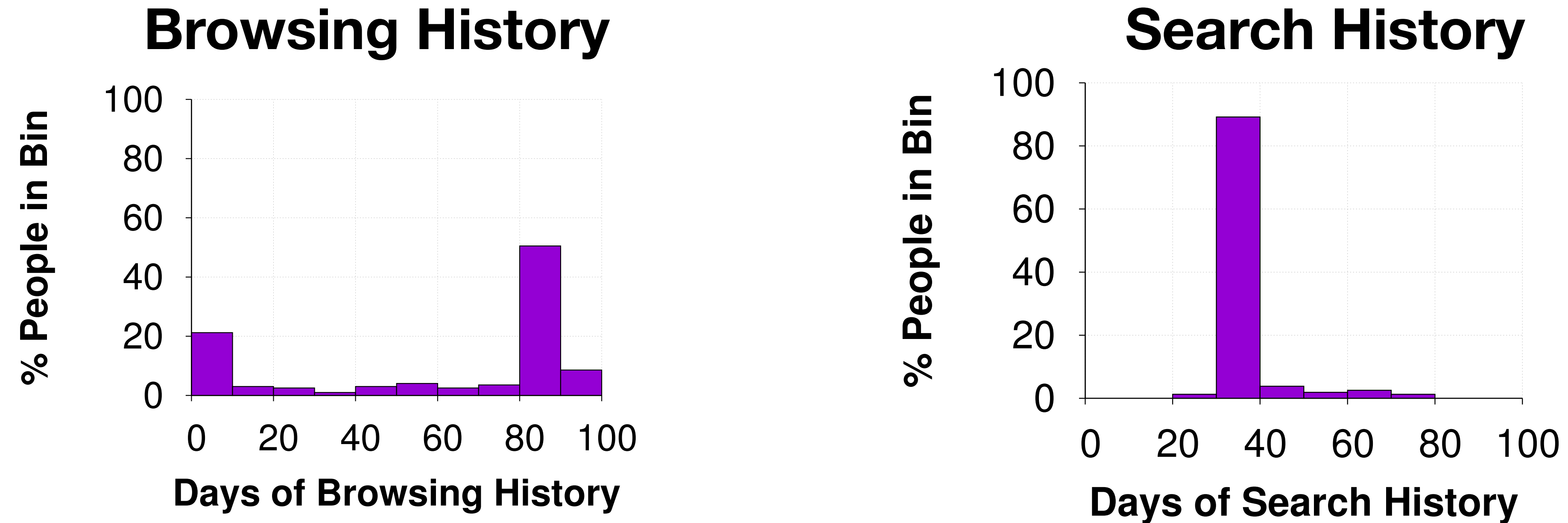


Fig: Amount of historical data collected from the participants

How Are The Inferences Drawn?

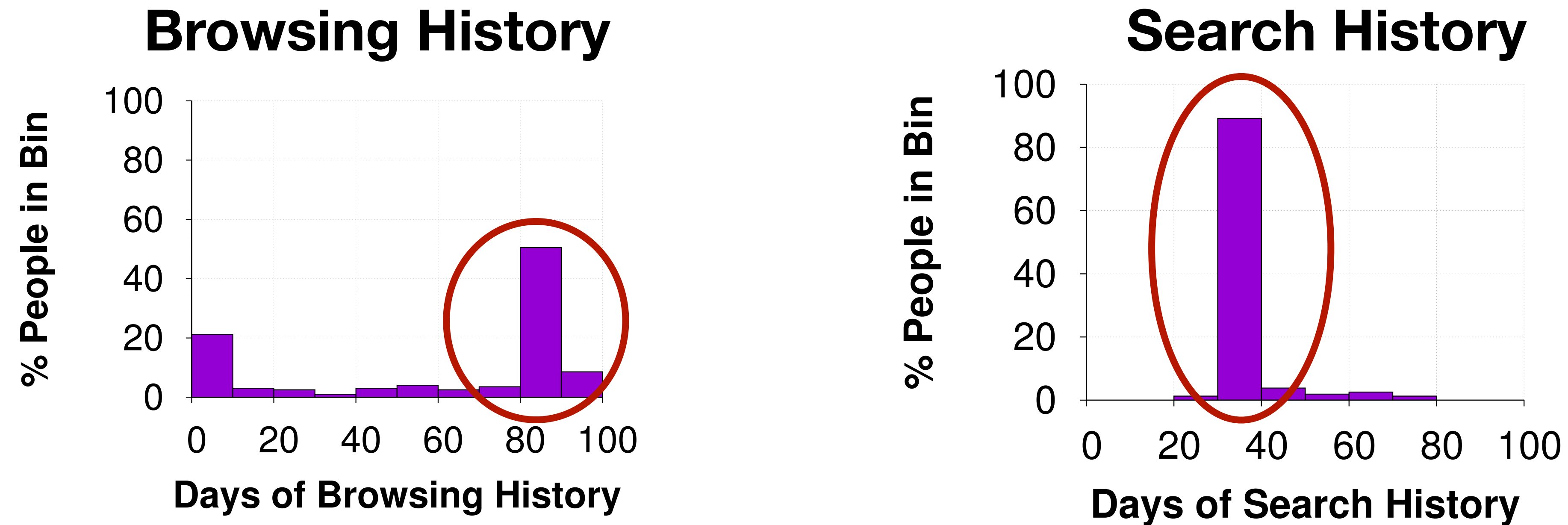


Fig: Amount of historical data collected from the participants

- 50% people had 80-90 days of browsing history
- 90% people had 30-40 days if search history

Domains From Browsing & Search History

Domains From Browsing & Search History

Browsing

- Out of 1.2M unique URLs, we extracted ~42K unique FQDNs

Domains From Browsing & Search History

Browsing

- Out of 1.2M unique URLs, we extracted ~42K unique FQDNs
- We used PhantomJS to collect trackers from these 42K FQDNs
 - We crawl home page + 5 additional pages

Domains From Browsing & Search History

Browsing

- Out of 1.2M unique URLs, we extracted ~42K unique FQDNs
- We used PhantomJS to collect trackers from these 42K FQDNs
 - We crawl home page + 5 additional pages
- Only considered those domains, where any of the APM trackers were present

Domains From Browsing & Search History

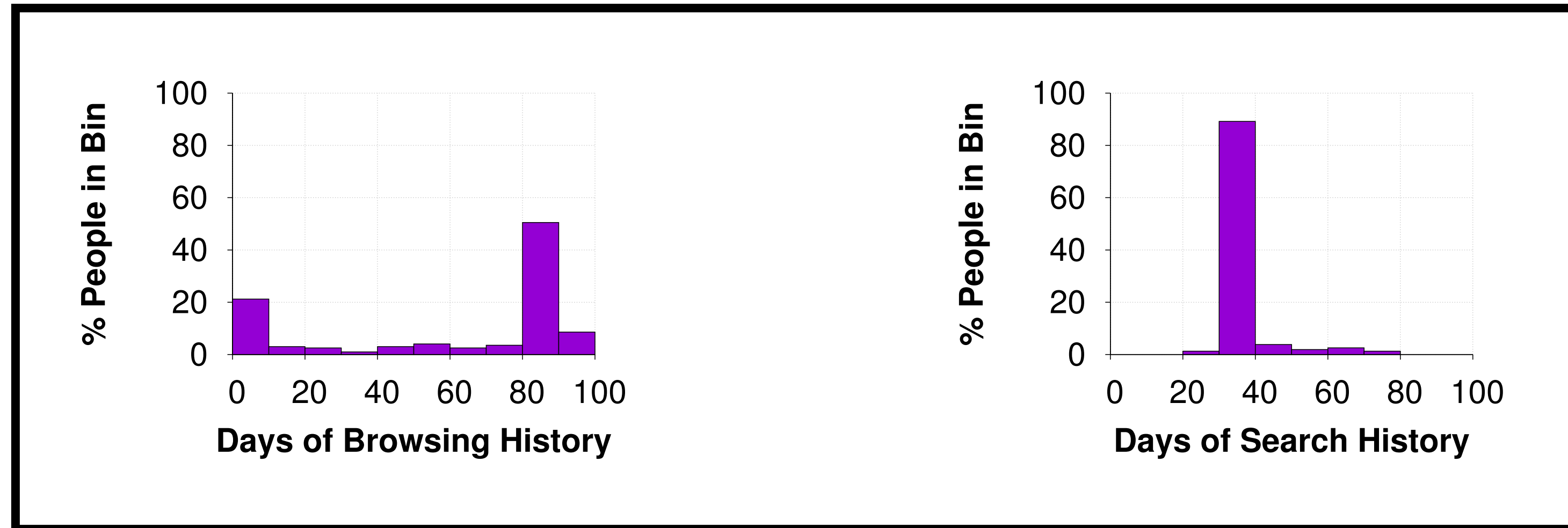
Browsing

- Out of 1.2M unique URLs, we extracted ~42K unique FQDNs
- We used PhantomJS to collect trackers from these 42K FQDNs
 - We crawl home page + 5 additional pages
- Only considered those domains, where any of the APM trackers were present

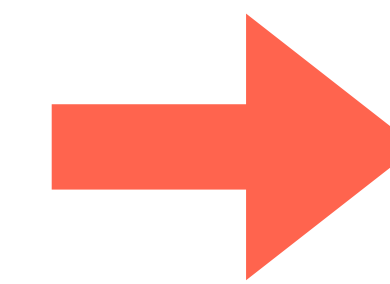
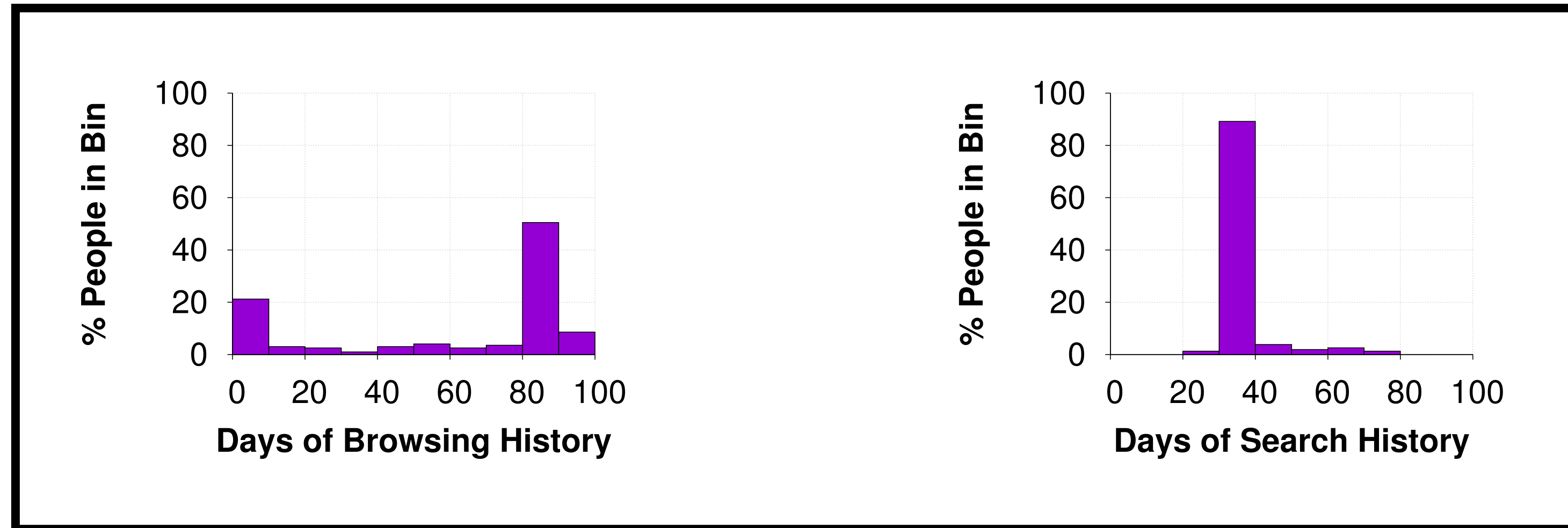
Search

- Considered the URL of the first search result

Domains Mapped to Common Space

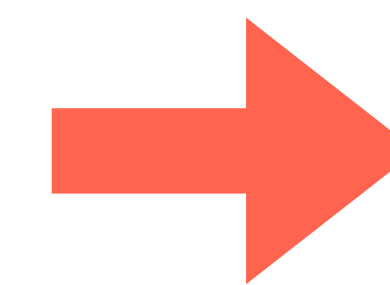
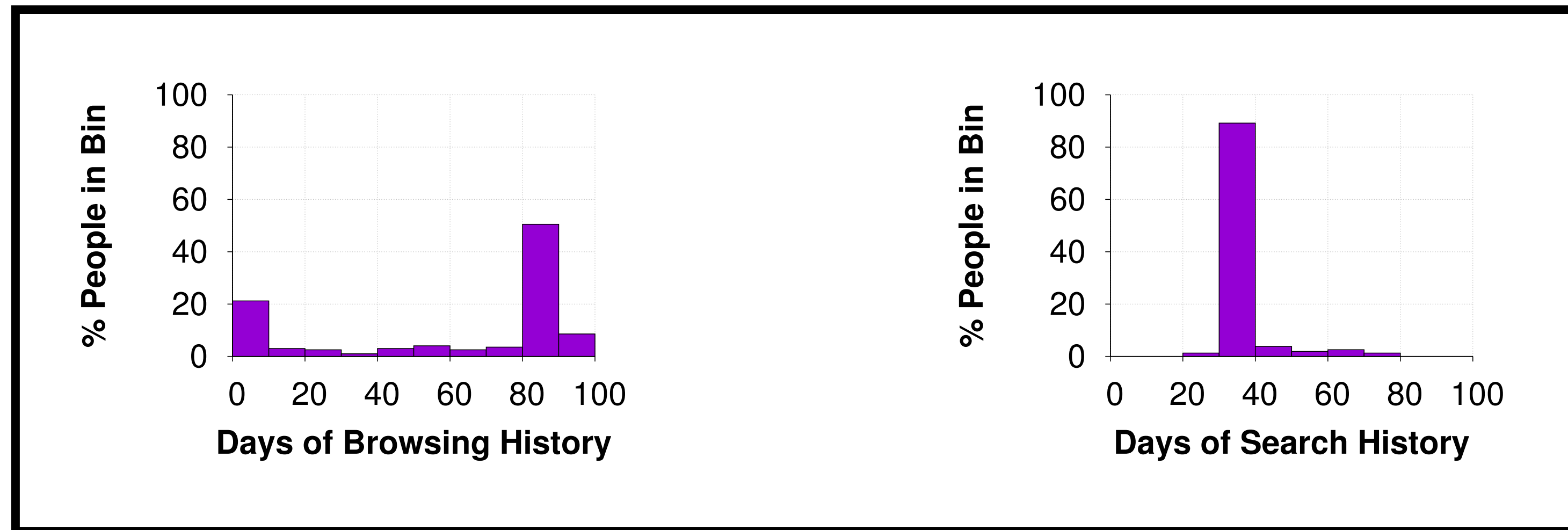


Domains Mapped to Common Space



51,500 unique domains

Domains Mapped to Common Space

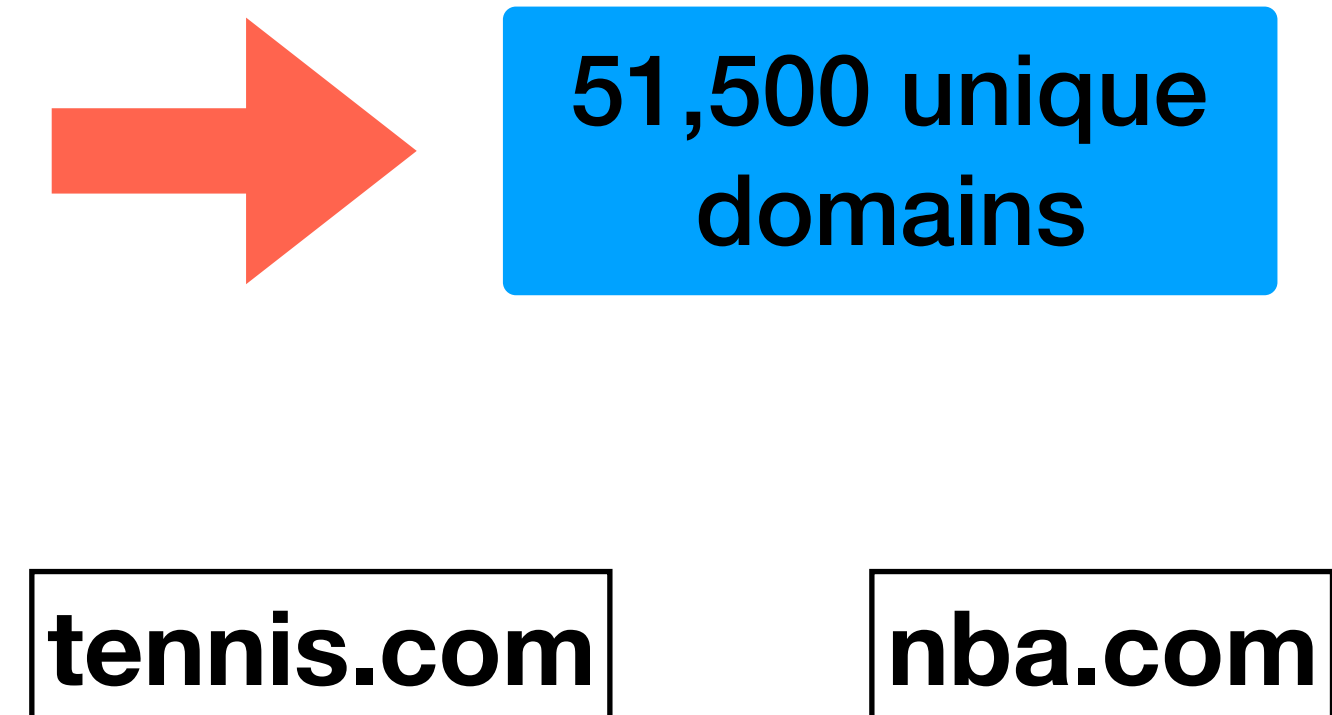
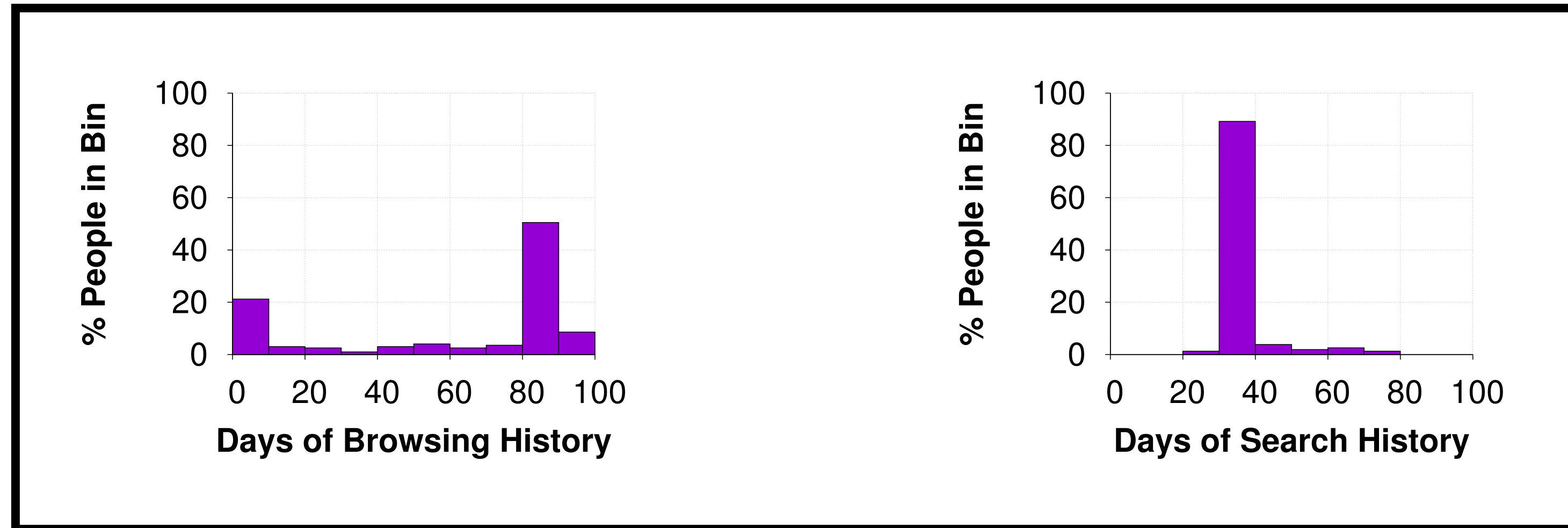


51,500 unique domains

We use SimilarWeb tool to map domains to (221) categories

- 77% success rate
- We then map each category to ODP category

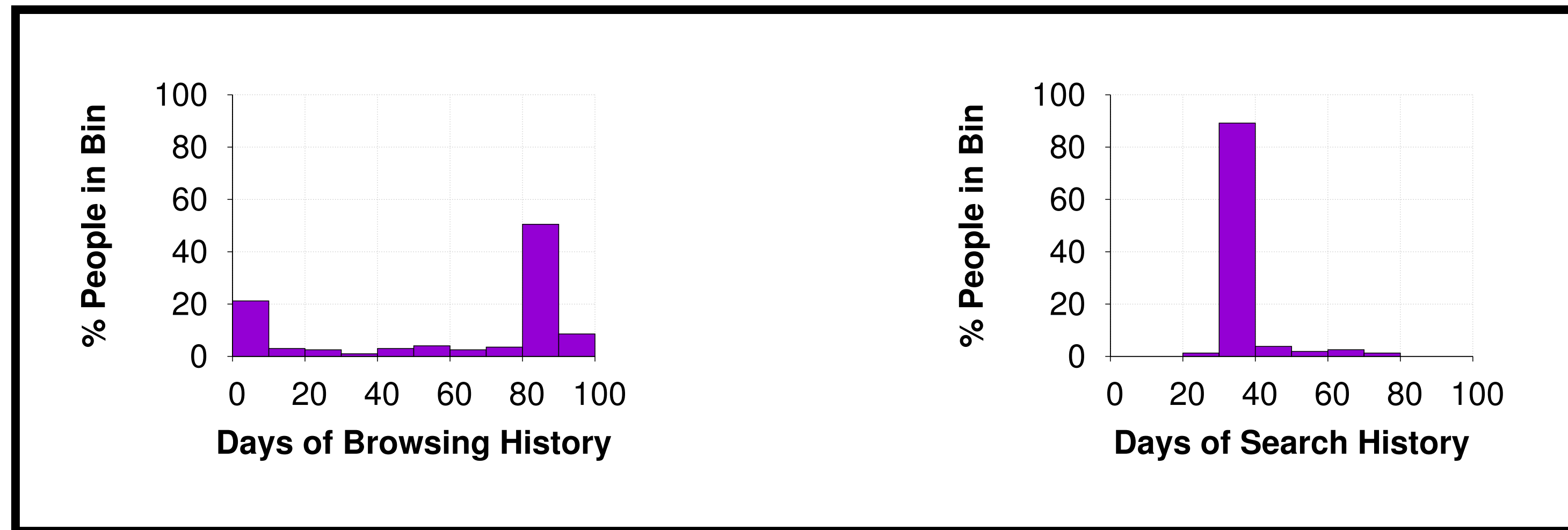
Domains Mapped to Common Space



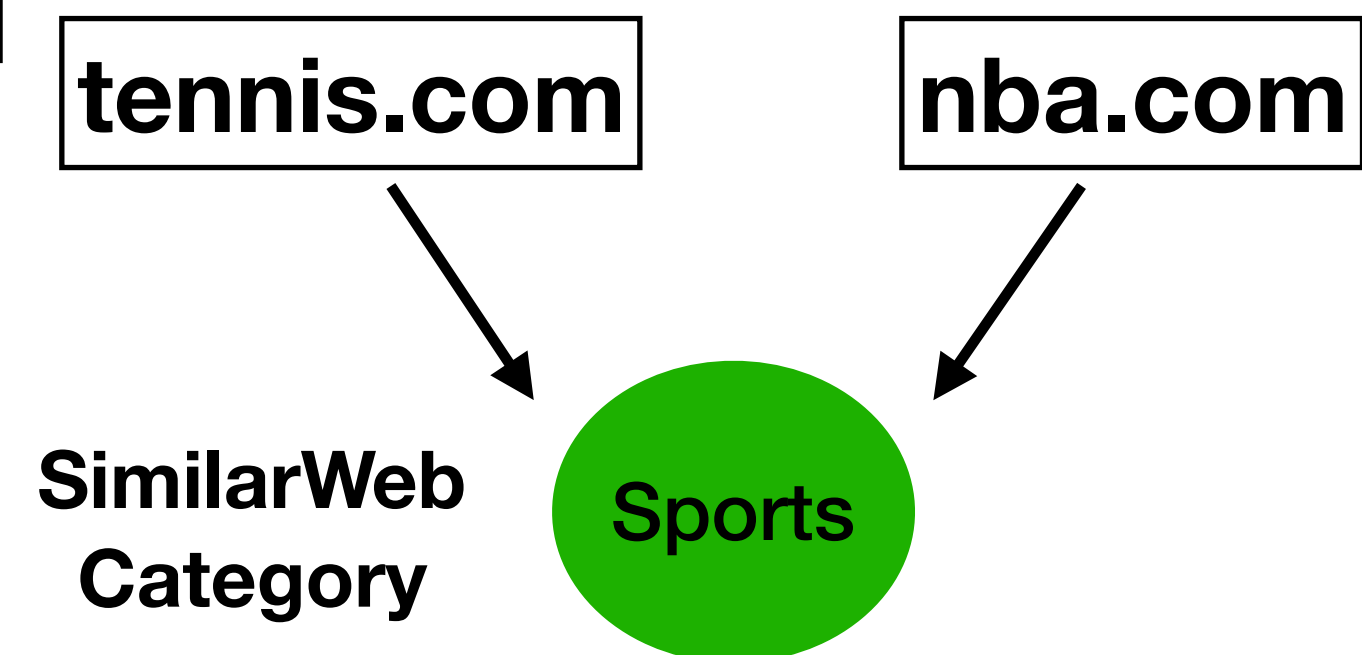
We use SimilarWeb tool to map domains to (221) categories

- 77% success rate
- We then map each category to ODP category

Domains Mapped to Common Space



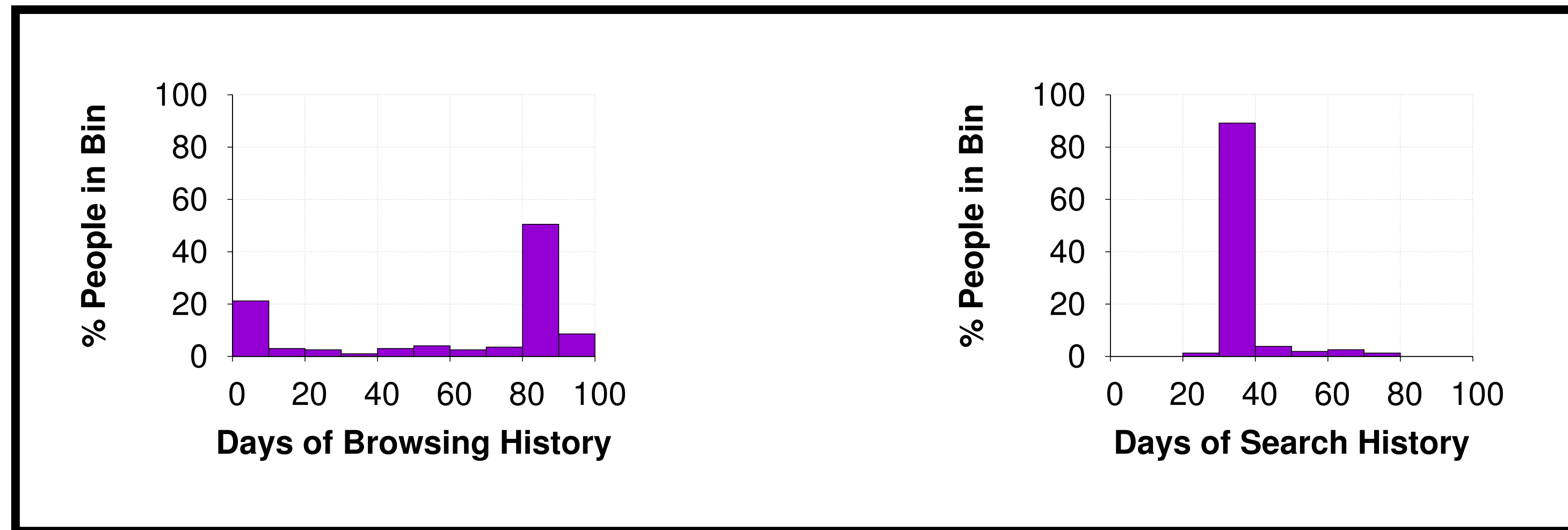
→ 51,500 unique domains



We use SimilarWeb tool to map domains to (221) categories

- 77% success rate
- We then map each category to ODP category

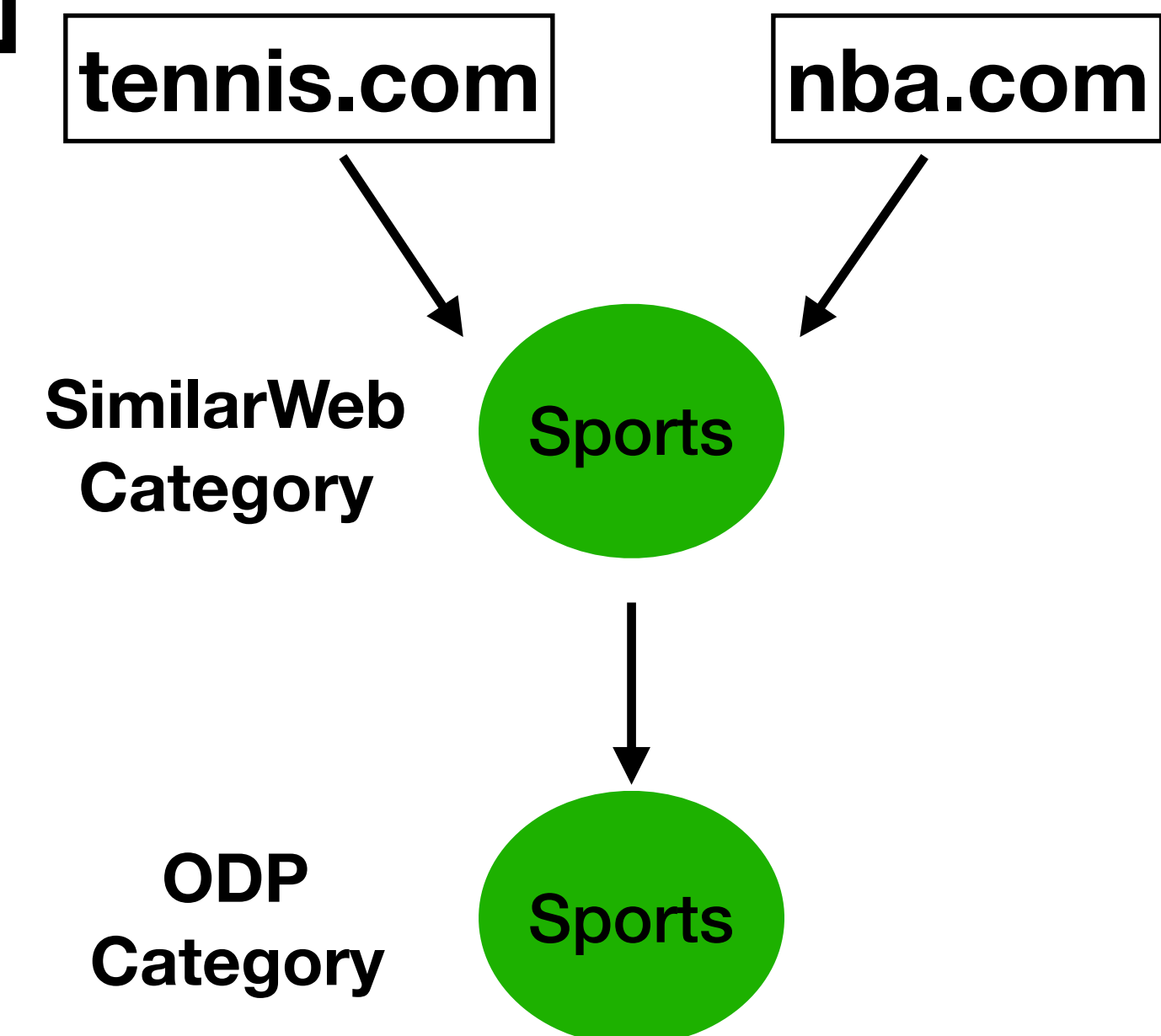
Domains Mapped to Common Space



→ 51,500 unique domains

We use SimilarWeb tool to map domains to (221) categories

- 77% success rate
- We then map each category to ODP category



Origins of Interests

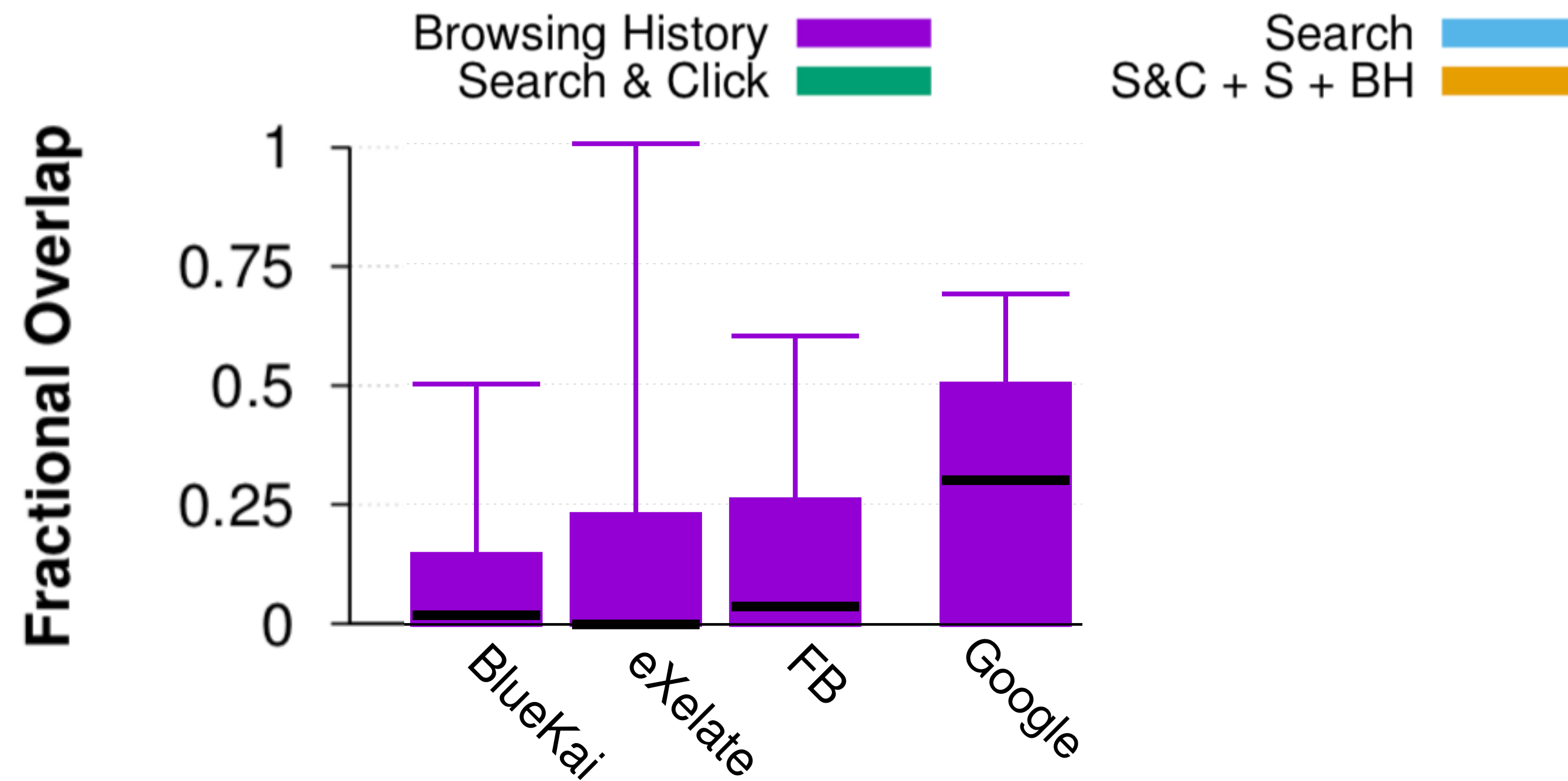


Fig: Overlap of Participants history with each APM (min, 5th, median, 95th, max)

Origins of Interests

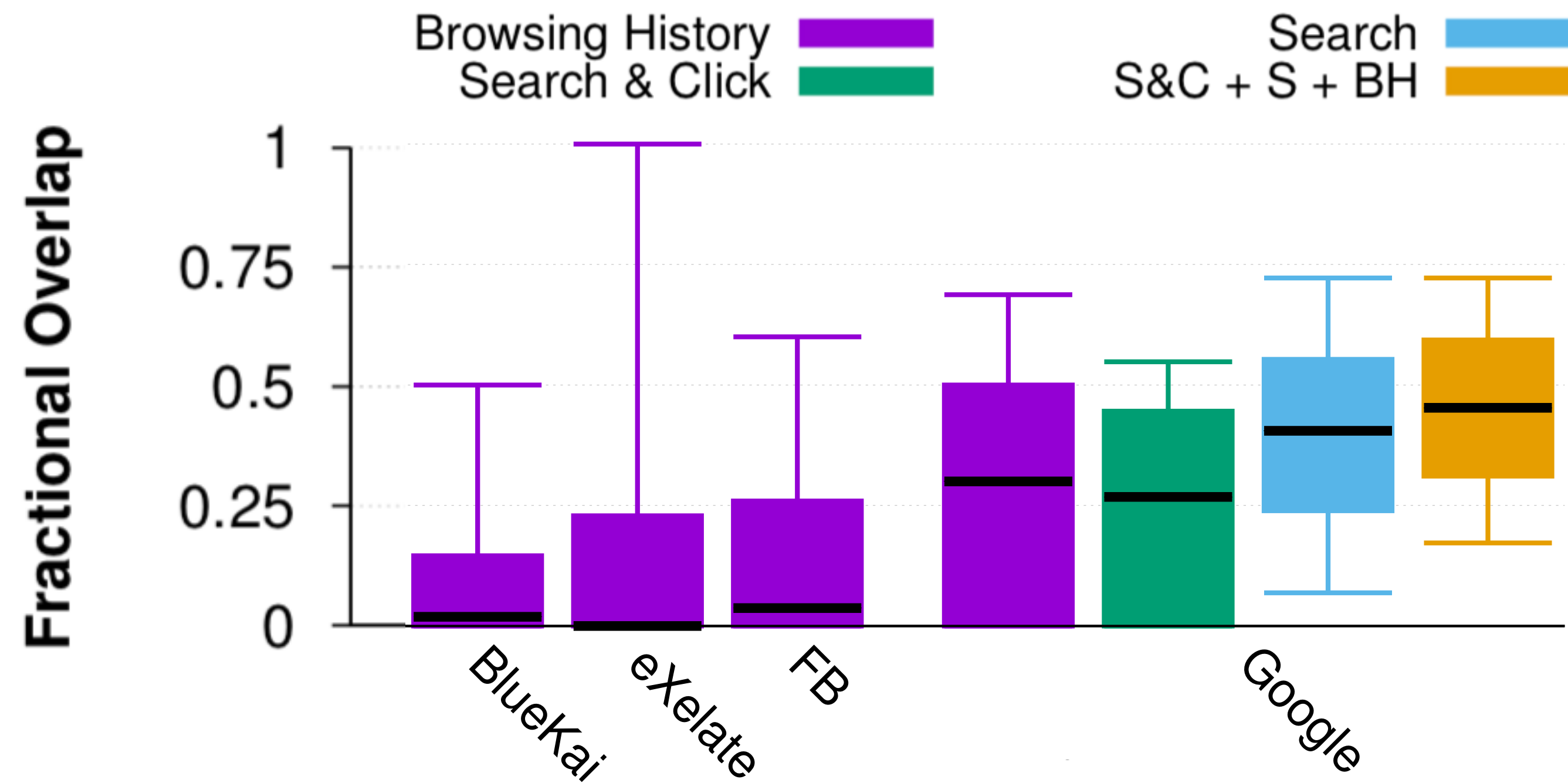


Fig: Overlap of Participants history with each APM (min, 5th, median, 95th, max)

Origins of Interests

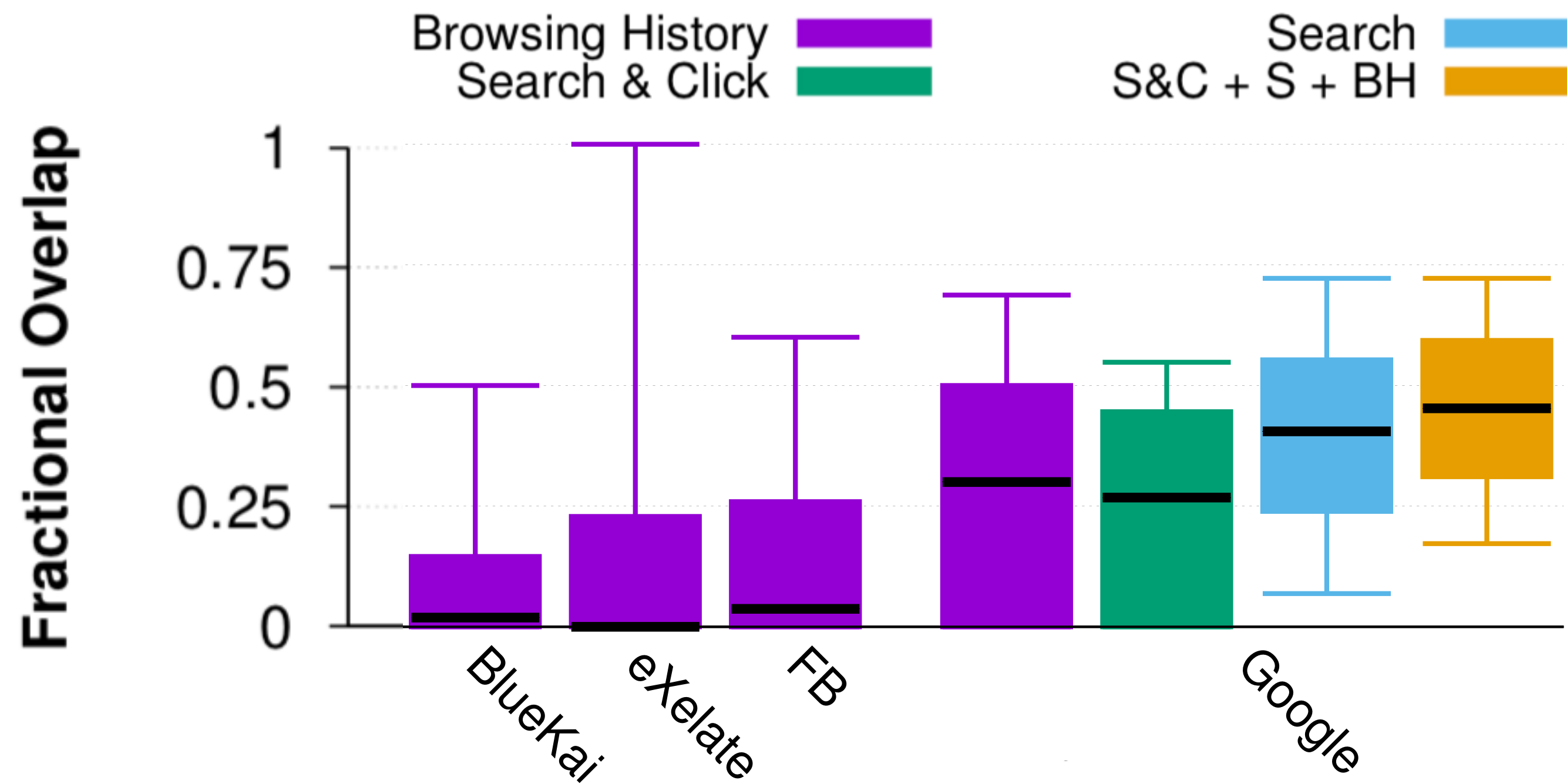
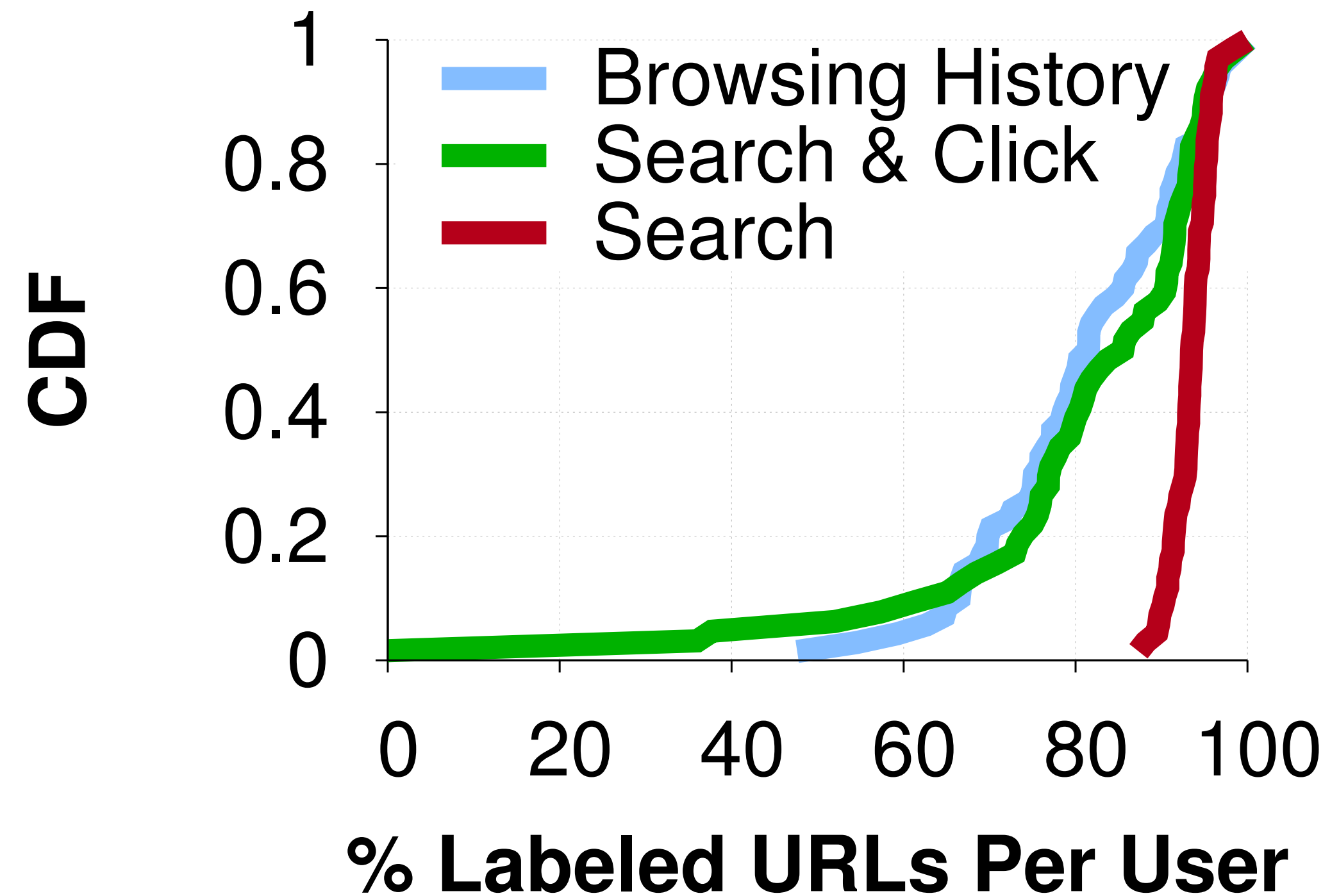
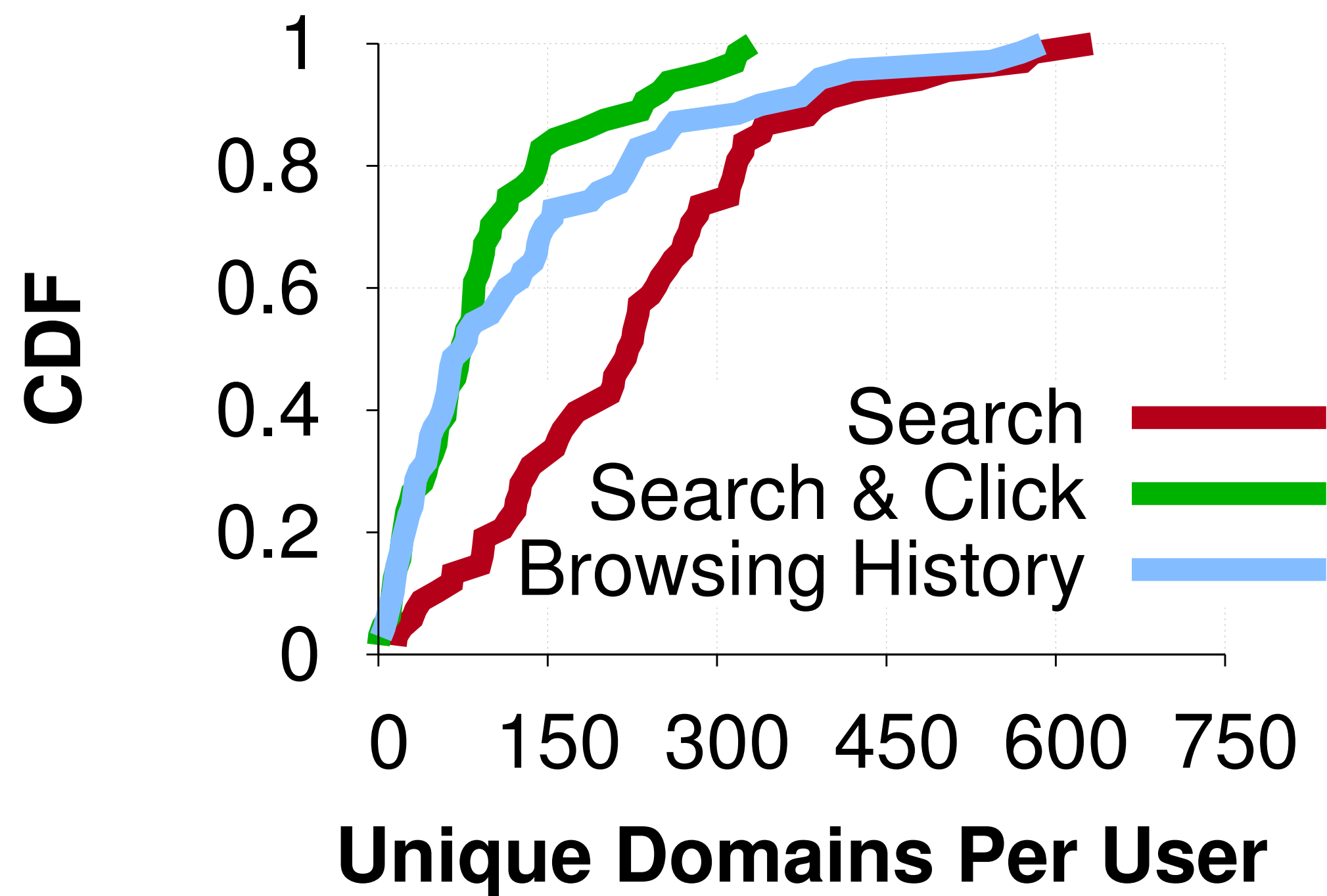


Fig: Overlap of Participants history with each APM (min, 5th, median, 95th, max)

Key Takeaways

- Browsing History explain <10% of interests, except for Google (30%)
- Search History does not add much to the explanation on top of BH

Browsing & Search History Domains



- More domains in Search as compared to Browsing
- Very high label rate for Search
- >75% Browsing domains labeled for 80% people

BlueKai Branded Data

alliant

acxiom

datalogix

acquireweb

lotame

affinity answers

experian

placeiq

adadvisor by neustar

tivo
